# CNS-BENCH: BENCHMARKING MODEL ROBUSTNESS UNDER CONTINUOUS NUISANCE SHIFTS

Anonymous authors

Paper under double-blind review

#### Abstract

One important challenge in evaluating the robustness of vision models is to control individual nuisance factors independently. While some simple synthetic corruptions are commonly applied to existing models, they do not fully capture all realistic distribution shifts of real-world images. Moreover, existing generative robustness benchmarks only perform manipulations on individual nuisance shifts in one step. We demonstrate the importance of gradual and continuous nuisance shifts, as they allow evaluating the sensitivity and failure points of vision models. In particular, we introduce CNS-Bench, a Continuous Nuisance Shift Benchmark for image classifier robustness. CNS-Bench allows generating a wide range of individual nuisance shifts in continuous severities by applying LoRA adapters to diffusion models. After accounting for unrealistic generated images through an improved filtering mechanism for such samples, we perform a comprehensive large-scale study to evaluate the robustness of classifiers under various nuisance shifts. Through carefully-designed comparisons and analyses, we find that model rankings can change for varying shifts and shift scales, which is not captured when averaging the performance over all severities. Additionally, evaluating the model performance on a continuous scale allows the identification of model failure points, providing a more nuanced understanding of model robustness. Overall, our work demonstrated the advantage of using generative models for benchmarking robustness across diverse and continuous real-world nuisance shifts in a controlled and scalable manner.

033

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

#### 1 INTRODUCTION

Machine learning models are typically validated and tested on fixed datasets under the assumption of
 independent and identically distributed samples. This, however, does not fully cover the true capabilities and potential vulnerabilities of models when deployed in dynamic real-world environments.
 The robustness in out-of-distribution (OOD) scenarios is important and decision-makers might need
 to know how models perform under various distribution shifts and severity levels in safety-critical
 scenarios. Therefore, it is crucial to continue building richer and more systematic benchmarks.

In the past few years, various benchmarks have been proposed to evaluate the robustness of computer vision models. One line of benchmarks manually collects data with nuisance shifts (Zhao et al., 2022; Hendrycks et al., 2021a; Wang et al., 2019; Geirhos et al., 2022; Barbu et al., 2019; Idrissi et al., 2022; Hendrycks et al., 2021b; Recht et al., 2019). Yet, such approaches are not scalable and often include only a small variety of nuisance shifts.

On the other hand, synthetic datasets offer opportunities to evaluate deep neural networks since various instances of an object class with specified context and nuisance shifts can be generated. While rendering pipelines allow precise control of several variables and are applied for benchmarking (Bordes et al., 2024; Shu et al., 2020; Kar et al., 2022; Li et al., 2023c), some nuisance shifts such as weather variations (*e.g.*, snow) are very hard to perform using traditional pipelines. While Hendrycks & Dietterich (2018) report accuracy drops for various types and levels of synthetic corruptions, they lack relevant real-world nuisance shifts.

Recent developments in diffusion models have enabled the application of generative models for training (He et al., 2022b; Fan et al., 2024) and benchmarking vision models (Mofayezi & Medghalchi, 2023; Metzen et al., 2023; Vendrow et al., 2023; Zhang et al., 2024). However, all

071



Figure 1: Benchmarking under continuous nuisance shifts. We evaluate the robustness of different models under gradually increasing nuisance shifts. This allows identifying the *failure point* (highlighted in red) of a model.

072 previous approaches define *categorical* or *binary* nuisance shifts by considering the existence or ab-073 sence of a shift, which contradicts their continuous realization in real-world scenarios. For example, 074 as shown in Fig. 1, the snow level in an environment can range from light snowfall to objects fully covered with snow. While one model might fail at all snow levels, a different model might only fail 075 when the object is heavily occluded. In most real-world applications, it is important to know the 076 expected performance at specific nuisance shift levels, rather than just a global accuracy drop. For 077 instance, an autonomous driving company may need to determine the fog density at which system performance falls below a critical threshold. Evaluating such failure points to probe the sensitivity 079 of models requires realizing continuous shifts. 080

081 To overcome this shortcoming, we establish a Continuous Nuisance Shift Benchmark for model robustness, dubbed as CNS-Bench. Specifically, we apply LoRA (Hu et al., 2021) adapters to 083 diffusion models to perform a continuous variation of specified nuisance shifts, and use them to benchmark a variety of classifiers along the following axes: (i) architecture, (ii) number of param-084 eters, (iii) pre-training paradigm and data. In contrast to previous works conducting analysis on 085 binary or categorical shifts, our study advocates multiple scales of shifts. We caveat that model 086 rankings can change when considering several scales. It is also essential to consider failure points, 087 *i.e.*, the shift severity at which a model fails. Thus, measuring robustness as a spectrum instead 880 of aggregating it into a single average metric allows a more comprehensive understanding of OOD 089 robustness (Drenkow et al., 2021; Hendrycks et al., 2021a). With our benchmark, we evaluate more 090 than 40 classifiers and demonstrate that a rigorously-designed generative benchmark allows system-091 atically studying the robustness behaviors of vision models in a controlled and scalable manner. 092

One essential requirement when using synthetic images for benchmarking is to ensure that the considered images correspond to the class distribution. Manually checking the quality of images to find those not aligned with the desired condition is still a common practice (Zhang et al., 2024). However, it has difficulty in scaling up the analysis (Hastie et al., 2009; Angelopoulos et al., 2023). Some approaches have been proposed for automatic filtering, but no standard datasets are available to evaluate filtering strategies. With this in mind, we also provide a dataset with manually annotated out-of-class (OOC) images. We show that our proposed filtering mechanism outperforms previous strategies in removing such problematic samples.

100 In summary, our work makes the following contributions: 1) We propose CNS-Bench to benchmark 101 vision models under continuous nuisance shifts. We publish a dataset with 14 diverse and realistic 102 nuisance shifts that represent various style and weather variations at five severity levels. In addition, 103 we also provide trained LoRA sliders for all shifts that can be used to compute shift levels in a fully 104 continuous manner. 2) We collect an annotated dataset to benchmark OOC filtering strategies and 105 propose a novel filtering mechanism that achieves higher filter accuracies than previous methods. 3) We evaluate the robustness of more than 40 classifiers along different axes and reveal multiple 106 valuable findings, underlining the importance of considering continuous shift severities of real-world 107 nuisance shifts.

## 108 2 RELATED WORK

109

**Robustness.** When referring to robustness, we consider the relative accuracy drop of a classifier *w.r.t.* interventions that alter images from a base distribution, building upon the formalism introduced in Drenkow et al. (2021). While the averaged accuracy drops provide an aggregated measure of the robustness, we consider the robustness *w.r.t.* specific nuisance shifts that can be modeled as causal interventions on the environment, the appearance, the object, or the renderer. We define such continuous interventions on metric scale.

115 116

Benchmarking robustness. Early approaches for benchmarking the performance and gener-117 alizability of models use fixed datasets, assuming independent and identically distributed sam-118 ples (Deng, 2012; Deng et al., 2009; Lin et al., 2014). However, this lacks scalability and fails to 119 capture the performance in real-world applications facing OOD scenarios. To tackle this challenge, a 120 line of research involves manually collecting data with nuisance shifts (Zhao et al., 2022; Hendrycks 121 et al., 2021a; Wang et al., 2019; Geirhos et al., 2022; Barbu et al., 2019; Idrissi et al., 2022; 122 Hendrycks et al., 2021b; Recht et al., 2019). However, these methods are often time-consuming 123 and labor-intensive since they require data crawling and human annotations. Moreover, they usu-124 ally capture only a subset of nuisance shifts that models may encounter in the real world and it is 125 challenging to ensure the disentanglement of these annotated nuisances.

126 Another line of research uses synthetic data for benchmarking, which offers the ability to generate 127 a large and diverse range of nuisance shifts with precise control (Hendrycks & Dietterich, 2018; 128 Bordes et al., 2024; Shu et al., 2020; Kar et al., 2022). However, these works are limited to nuisances 129 that can be easily modelled (e.g., lighting, fog, occlusions) or restricted to what can be expressed in 130 rendering pipelines. Recent developments in diffusion models shed light on creating realistic and 131 diverse synthetic benchmark datasets (Mofayezi & Medghalchi, 2023; Metzen et al., 2023; Vendrow 132 et al., 2023; Zhang et al., 2024) with realistic data and more possibilities to control nuisances (e.g., 133 text-guided corruptions, counterfactual). In our work, we propose a framework to benchmark vision models w.r.t. nuisance shifts under multiple severity levels. To address the need to remove OOC 134 images from generative models, which are essential for benchmarking applications, we additionally 135 propose a novel strategy to remove such samples from the dataset. 136

137 138

139

## 3 CONTINUOUS NUISANCE SHIFT BENCHMARK

In this section, we present how CNS-Bench is created. We first discuss the strategy to replicate
 the in-domain distribution in Section 3.1. We then present our methodology to perform continuous
 shifts to evaluate the model's sensitivity to various nuisance factors in Section 3.2. Finally, we detail
 our filtering dataset and the selected filtering strategy in Section 3.3.

144 145

146

#### 3.1 REPLICATING THE IMAGENET DISTRIBUTION

We aim to evaluate a model's robustness to specific nuisance shifts that alter the base ImageNet 147 (Deng et al., 2009) distribution  $p(X_{IN}|c)$ , conditioned on an ImageNet class c. However, as pointed 148 out by Vendrow et al. (2023), the distribution of Stable Diffusion (SD) (Rombach et al., 2022) gen-149 erated images  $p(X_{SD}|c)$  differs from the ImageNet distribution, significantly lowering classification 150 accuracies. To generate images that are more similar to the ImageNet images, we apply textual inver-151 sion (Gal et al., 2023) to learn new "words" in the embeddings space of a text encoder that capture 152 the ImageNet-specific class concepts. Specifically, these text embeddings are optimized by mini-153 mizing the noise prediction error of diffusion models  $||\epsilon - \epsilon_{\psi}(\cdot, f_{\psi}(c))|^2$  with the text encoder  $f_{\psi}(\cdot)$ 154 and parameters  $\psi$  for all diffusion time steps. We call this distribution IN\*:  $p(X|c) = p(X_{IN*}|c)$ .

155 156

157

#### 3.2 CONTINUOUS NUISANCE SHIFTS FOR BENCHMARKING

To evaluate the robustness of vision models *w.r.t.* continuous nuisance shifts, the following characteristics are desirable: (i) The severity of the considered shift can be controlled, allowing the estimation of the shift scale where a considered model fails. (ii) Realizing a nuisance shift should not come along with variations that might alter the class identity. (iii) The variations should be subtle, allowing a fine-grained analysis also for specific images.



Figure 2: Qualitative examples for prompt-based and LoRA-based shifts with out-of-class samples. On the left, we present two images in a different a) style and b) weather condition generated from a text prompt a) "fox in cartoon style" and b) "birdhouse in heavy snow", respectively. On the right, we show the gradual variation performed by our LoRA sliders. a) Unlike the prompt-based shift, our LoRA sliders successfully generated images showing a gradual shift. b) Our LoRA sliders sometimes result in out-of-class (OOC) samples for higher scales, as depicted with the orange box.

175

176

177

178

179

Realizing continuous nuisance shifts. A natural way to perform synthetic nuisance shifts are methods based on text prompts (Metzen et al., 2023; Liu et al., 2023; Vendrow et al., 2023). They follow the two prompt (2P) templates: "A picture of a <class>" and "A picture of a <class> in <shift>". However, this approach does not allow the gradual increase of a nuisance for a given image. In addition, the generated shifts largely vary for different seeds and classes when applying the prompt addition "in <shift>"—for some seeds, the generated shift is more prominent, while for others, it is barely visible. Additionally, the semantic structure of the generated image can be significantly changed.

We leverage LoRA (Hu et al., 2021) adapters that represent low-rank matrices added to the original 190 weight matrices to perform continuous shifts. Such adapters are trained to characterize the effect of 191 a considered nuisance shift. Gandikota et al. (2023) propose a strategy to learn concept sliders using 192 LoRA adapters that allow a continuous modulation of the considered concept, which is achieved by 193 learning low-rank matrices that increase the expression of a specific attribute when applied to a class 194 concept c. The low-rank parameters  $\theta_{\text{LoRA}}$  modify the original model parameters  $\theta$  to  $\theta^* = \theta + s$ . 195  $\theta_{\text{LoRA}}$  with scale s are trained to capture a concept of interest  $c_+$ :  $P_{\theta^*}(X|c) \leftarrow P_{\theta}(X|c) \cdot P_{\theta}(X|c_+)^{\eta}$ , 196 where  $\eta$  refers to weighting factor that is fixed during training. Following Gandikota et al. (2023), we 197 optimize with the MSE objective (Sohl-Dickstein et al., 2015) using the Tweedie's formula (Efron, 2011) and the reparametrization trick (Ho et al., 2020) by formulating the scores as a denoising 199 prediction  $\epsilon(X, c, t)$  with the diffusion timestep t:  $MSE(\epsilon_{\theta^*}(X, c, t); \epsilon_{\theta}(X, c, t) + \epsilon_{\theta}(X, c_+, t))$ . We model the class concept c and the nuisance concept  $c_+$  by two text embeddings "<class>" 200 and "<class> in <shift>". Different to (Gandikota et al., 2023), we specifically use class 201 concepts c that are acquired from the IN\* distribution. After training, the learned LoRA adapters 202 capture the direction between the two language concepts, i.e., they characterize attributes of the 203 concept of interest  $c_+$ . Weighting their effect using the scale s modulates the effect of the applied 204 shift. Gandikota et al. (2023) stated that the LoRA adapters generalize to other concepts and images. 205 We found that learning class-specific LoRA sliders produces higher-quality shifts. This choice also 206 allows capturing the class-specific characteristics and confounders of the considered shifts that occur 207 in the real world. Hence, we train separate LoRA adapters for each ImageNet class and shift. As 208 qualitatively shown in Fig. 2, applying these learned directions enables gradual nuisance shifts. We 209 show examples of more shifts in Fig. 33 and Fig. 34.

Following Mokady et al. (2023); Gandikota et al. (2023), we evaluate the shift severity based on the CLIP similarity of the generated image to the text prompt describing the shift, *i.e.*, "A picture in <shift>". Similarly, we also compute the CLIP (Radford et al., 2021) similarity to the class prompt "A picture of a <class>". To measure the performed shift, we compute the CLIP shift difference by  $\Delta$ CLIP<sub>shift</sub>( $I_k$ ,  $I_0$ ) =  $\cos$  (CLIP<sub>img</sub>( $I_k$ ), CLIP<sub>text</sub>("in {shift}")) -  $\cos$  (CLIP<sub>img</sub>( $I_0$ ), CLIP<sub>text</sub>("in {shift}")) for the generated image with scale 0 and scale k, and similarly for the class similarity. In contrast to simply applying a second text prompt (2P) to
perform a *binary* shift, our LoRA adapters allow performing a
variety of shift scales, as measured by the CLIP shift difference
(Section 3.2). This allows gradual shifts (Fig. 2).

220 Activating the LoRA adapter at different time steps throughout 221 the diffusion process will modulate the effect of the adapter 222 on the generation process (Meng et al., 2021). If the LoRA 223 adapter is active for all noise steps, it will significantly in-224 fluence the semantic structure and the appearance of the gen-225 erated image, while deactivating the adapter for earlier time 226 steps will keep the semantic structure. Since we aim to perform more fine-grained edits that do not heavily change the 227 semantic structure, we deactivate the LoRA adapter for early 228 steps. This allows applying edits where the semantic structure 229 remains similar but the appearance changes (e.g. Fig. 2a). 230



Figure 3: Average  $\triangle$  CLIP evaluation for the snow shift. Our sliders perform a gradual shift, while a naive application (2P) only allows *binary* shifts.

Since our sliders do not explicitly exclude confounding variables, the applied shifts may also affect
confounders inherently present due to biases in the training data. For example, as shown in Fig. 34c,
using the *in dust* slider unintentionally removes half of the people, and in Fig. 33c, the background
no longer represents a forest. Consequently, failures in our subsequent analysis cannot always be
solely attributed to the nuisance concept itself.

Failure point concept. We define a failure point  $s = \min\{S \in \mathbb{R} | f(X(S)) \neq c\}$  as the smallest shift scale where a classifier f(X(s)) fails to correctly classify an image X(s) with a class c and a scale s of a considered shift. The failure point distribution captures the ratio of failed samples for the considered scales. We estimate this distribution in our work with a histogram, where the number of elements in one bin  $I_k$  is computed by  $H(I_k) = \sum_{n=1}^N \mathbb{1}_{I_k}(s_n)$  with the indicator function  $\mathbb{1}(\cdot)$ and the scale of the *n*-th element of the set of images with N images. We compute and report the ratio of failure points for each scale s, dividing  $H(I_k)$  by the number of considered images N.

243 244

245

3.3 FILTERING DATASET AND STRATEGY

To evaluate filtering strategies for removing out-of-class (OOC) samples, we collect a manually labeled dataset. This section presents this dataset and the selected filtering strategy.

248 Filtering of OOC samples. Current diffusion models allow the generation of diverse and realistic 249 images  $x \sim p(X|\mathbf{z})$  that are conditioned on  $\mathbf{z} = [c, s_i]$ , which involves the considered ImageNet 250 class  $c \in \mathbb{N} \mid 1 \le c \le 1000$  and the variable  $s_i \in \mathbb{R}$  corresponding to the severity of a considered 251 nuisance shift i. However, due to their probabilistic formulation, the generated sample might deviate 252 from the condition z. For benchmarking applications, we are particularly concerned about gener-253 ated samples deviating from the original class c, *i.e.*, the considered class cannot be characterized 254 anymore (c.f., Fig. 2). We call such samples "OOC" samples (Metzen et al., 2023). Evaluating the 255 sensitivity to specific nuisance shifts requires removing the OOC samples generated by the shift's application. Therefore, we collect a dataset of generated images to evaluate the sliding process and 256 strategies to automatically remove OOC samples. 257

258 Dataset for evaluating OOC filtering strategies. To evaluate various OOC filtering strategies, we 259 manually label a dataset consisting of 18k generated images with two shifts, five scales, and 100 260 random ImageNet classes. We select *snow* as one weather variation and *cartoon* as one style shift to 261 represent two rather different nuisance shifts. Before manually labeling the dataset, we remove easy samples that have a high CLIP text alignment and are classified correctly by multiple classifiers. 262 Then, all hard images are labeled by two human annotators, where each annotator can choose from 263 the following labels: "class", "partial class properties", and "not class". More details on the labeling 264 strategy and the dataset statistics are provided in Appendix A.6. 265

OOC filtering strategy. A filter serves its purpose if it removes all OOC samples, corresponding to a high true positive rate (TPR), while not removing too many in-class samples, corresponding to a low false positive rate (FPR). Instead of simply applying a CLIP threshold as in Vendrow et al. (2023), we consider a combinatorial selection approach, which requires two out of four filters to be active. For the first and the second filter, we consider text alignment to "A picture of a <class>"



279

Figure 4: Accuracy drops and failure point ratios of a ResNet-50 classifier on OOD-CV and our benchmark. *Left*: Accuracy drops on OOD-CV and various scales of our benchmark (in the value range [0,1]). Horizontal lines show the average score for each weather nuisance of the OOD-CV test dataset, while our benchmark allows identfying the performance drop at various shift scales. *Right*: Distribution of failure points. Our continuous nuisance shifts allow identifying the scales that result in a failure and models fail earlier for fog, potentially due to heavier occlusions than snow and rain.

and "A picture of a <class> in <shift>", respectively, computed via CLIP. For the third and fourth filter, we measure the cosine similarity to the starting images using the CLIP image encoder and the class tokens of DINOv2 (Oquab et al., 2023), respectively. We select the filtering threshold for each filter such that 90% of the labeled OOC samples are removed. Note that none of these filters are trained on ImageNet data.

291 292

293 294

295

296

297

299

#### 4 EXPERIMENTS

In this section, we discuss our benchmark results. First, we compare our bechmarking strategy with the OOD-CV benchmark. Then, we perform a large-scale analysis by evaluating more than 40 ImageNet classifiers on CNS-Bench.

## 4.1 Comparing Continuous Shifts with OOD-CV Dataset

300 Zhao et al. (2022; 2024) introduce OOD-CV to measure out-of-distribution (OOD) robustness, a 301 benchmark dataset that includes OOD examples of ten object categories for five different individual nuisance factors (e.g., weather) on real data. OOD-CV is the only real-world dataset that provides 302 accurate labels of various individual weather shifts. This allows comparing our generated images 303 with real-world weather realizations of the considered shifts. We use our trained LoRA adapters 304 to create a benchmark for the OOD-CV classes and scales up to 3.0 to directly compare with the 305 original manually labeled dataset. We refer to the supplementary for exemplary images of both 306 benchmarks and CLIP alignments to the considered shifts. 307

First, we train a ResNet-50 classifier on the training set of the OOD-CV benchmark. Then, we 308 evaluate the performance on our data and the OOD-CV benchmark. Fig. 4 presents the results for 309 each nuisance independently. The accuracies remain more or less constant with an accuracy around 310 95% up to a nuisance scale of 1.5. From a nuisance scale of 2.0, the accuracy starts dropping, with 311 the nuisance of fog having the biggest impact. This could be explained by the fact that fog can lead to 312 severe occlusion, while rain and snow can be considered as corruption factors. We hypothesize that 313 the partially bigger drop for the OOD-CV benchmark is due to a major limitation of its dataset: The 314 nuisances are not completely disentangled, and part of the accuracy drop originates from various 315 other factors (e.g., image quality, image size, and noise). In contrast, our benchmark allows for 316 fine-grained control of nuisances with multiple shift levels, leading to a more complete and scalable 317 analysis of the model's performance.

318

# 4.2 EVALUATED MODELS AND EXPERIMENTAL SETUP 320

We use our large-scale benchmark to evaluate the models along the following axes:

(i) Architecture. To compare architectures with a comparable number of parameters, we consider
 both CNN and ViT architectures with different training recipes: ResNet-152 (He et al., 2016), ViT-B/16 (Dosovitskiy et al., 2020), DeiT-B/16 (Touvron et al., 2021), DeiT-3-B/16 (Touvron et al.,



Figure 5: Classification accuracy drops on the labeled and filtered datasets. The accuracy drop curves of ResNet-50 and DINOv2 classifiers on the filtered and the labeled dataset are comparable, demonstrating the effectiveness of our automatic filtering strategy. The ratio of misclassified images is given in the value range [0,1]. We provide results for more classifiers in Fig. 8.

332

333

334

335

2022), and ConvNeXt-B (Liu et al., 2022). All models are trained in a supervised manner.

(ii) *Model size*. For ViT, we consider the small, medium, base, large, and huge variants of DeiT-3.
For CNN, we consider the ResNet variants: 18, 34, 50, 101, and 152.

(iii) *Pre-training paradigm and data.* We evaluate a set of models with the same backbone but different pre-training strategies. The following models are pre-trained on IN1k with a self-supervised objective: MAE (He et al., 2022a), DINOv1 (Caron et al., 2021), and MoCov3 (Chen et al., 2021). To study the impact of more data during training, we compare their performance to a supervised model that is trained only on ImageNet-1k and a supervised model that is pre-trained on ImageNet-21k. All Transformer-based models use ViT-B/16 as the backbone. Furthermore, we evaluate an ImageNet-trained diffusion classifier (Li et al., 2023b) on a smaller subset due to its heavy computational cost.

- Metrics. We typically report the average accuracy drops, i.e., the ratio of failed images, averaged over the images of one shift or all shifts in the value range [0, 1]. In Table 1, we report the mean relative corruption error (rCE) as introduced by Hendrycks & Dietterich (2018). It is defined by the average over all relative corruption errors  $CE_{shift} = \left(\sum_{s} E_{shift,s}^{f} - E_{shift,0}^{f}\right) / \left(\sum_{s} E_{shift,s}^{alex} - E_{shift,0}^{alex}\right)$
- with the average error E for scale s, model f, and a specific shift.
- 355 **Slider details.** As pointed out in Section 3.1, we use textual inversions to replicate the ImageNet 356 distribution. To evaluate the relevance of this approach, we generate 200 images of 100 randomly 357 selected ImageNet classes using standard SD2.0 and SD2.0 with the textual inversions of IN\*. To 358 illustrate the distribution gap, we compute the accuracies for ResNet-50 (DeiT). They achieve an accuracy of 68.2% (71.6%) for the SD distribution and 74.1% (79.1%) for the IN\* distribution, 359 which equals accuracy drops of 6% (8%) for both classifiers. This is significantly closer to the 360 performance on the original ImageNet distribution. We perform all the following experiments using 361 the IN\* distribution. We use SD2.0 and we activate the LoRA adapters with the selected scale for 362 the last 75% of the noise steps. 363
- Due to the computational complexity, we perform sliding for 100 classes. To get an estimate of the robustness on the full scale of ImageNet, we classify based on 1000 classes using off-the-shelf classifiers without applying classifier masking, as done by Hendrycks et al. (2021a). We ablate how the number of classes influences the robustness evaluations in Appendix A.5.2.
- The selection of the shifts is mainly inspired by ImageNet-R Hendrycks et al. (2021a) (8 shifts) and the OOD-CV dataset Zhao et al. (2022) (6 shifts) to consider a diverse set of nuisance shifts that modulate the appearance and style or the background and occlusion. Specifically, we consider the following 14 shifts: cartoon style, plush toy style, pencil sketch style, painting style, design of sculpture, graffiti style, video game renditions style, style of a tattoo, heavy snow, heavy rain, heavy fog, heavy smog, heavy dust, and heavy sandstorm.
- Filtering details. Our OOC filtering mechanism reaches a TPR of 87.9% and an FPR of 12.0% with an accuracy of 88.0%, while the naive CLIP-based thresholding reaches a TPR of 89.9% and an FPR of 35.7% with an accuracy of 65.1%. We plot the classification accuracy of DINOv2-R and ResNet-50 for the labeled and the filtered versions in Fig. 5. We observe comparable accuracy drops on both the manually-labeled and the filtered datasets. To further support the realism of our generated

Table 1: rCE along the model axes. We choose the average relative corruption error Hendrycks
& Dietterich (2018) as a single metric to measure the performance of a model on our benchmark
(lower is better). We provide results for all models in Table 2.

Architec	ture	Size	e	Pre-Training		
ConvNext	0.686	DeiT3-S	0.747	DINOv1-IN1k	0.636	
DeiT3	0.610	DeiT3-M	0.758	MAE-IN1k	0.732	
DeiT	0.746	DeiT3-B	0.610	MoCov3-IN1k	0.669	
RN152	0.790	DeiT3-L	0.574	SUP-IN1k	0.926	
ViT	0.926	DeiT3-H	0.583	SUP-IN21k-1k	0.722	



(a) Accuracy drops averaged over the whole benchmark. Architecture (*left*): We show models with the same training data and similar parameter counts. The selection of the architecture influences the accuracy drop. Model size (*center*): We show DeiT3 with various numbers of parameters. Increasing the model capacity results in lower accuracy drops. Pre-training paradigm and data (*right*): We show different pre-training paradigms: supervised, self-supervised (MAE, DINO, MoCo), and more data (IN21k), all using ViT-B/16. We present results for all shifts in Fig. 9.



(b) Accuracy drops for three selected shifts. Models exhibit varying performance changes depending on the considered shifts. For snow and painting shifts, the ranking of the models changes. In contrast, the cartoon style shift results in a consistent model ranking. However, the OOD performance on cartoon-shifted images is drastically worse than the other shifts.



(c) Ratio of failure points per scale for various models and shifts. The distribution allows inferring at which scales various models fail most often. Different models fail at varying stages depending on the considered shifts. While the number of failure points gradually increases for the snow shift, most failure points occur around scale 1.5 for the cartoon style shift. We present results for all shifts in Fig. 10.

Figure 6: **Evaluation of accuracy drops and failure points.** We plot the averaged accuracy drops and failure points of selected models and provide the results of all evaluated models in Appendix A.2.

431 images, we fine-tune ResNet-50 with our data and show more than 10% gains on ImageNet-R (see Appendix A.3).



Figure 7: **Relation between ID and OOD accuracy.** We report the slope of the linear fit between ID and OOD accuracy using 16 supervised ImageNet-trained models for all evaluated shifts. The relation varies for different shifts and scales between 0.5 and 2.5.

446 4.3 ANALYSIS AND FINDINGS

443

444

445

447

In this subsection, we discuss the main findings on our benchmark. Following Hendrycks et al. (2021a), we report the accuracy drops for 5 scales averaged over 14 diverse shifts as a measure of robustness in Fig. 6a. Table 1 compares models using the average relative corruption errors as proposed by Hendrycks & Dietterich (2018). We also provide results for three exemplary shifts in Fig. 6b. In addition, we report the distribution of failure points in Fig. 6c. We provide more evaluations in Appendix A.2.

454 Considering multiple scales of a shift allows a more nuanced analysis of OOD robustness. We 455 present the accuracy drops for multiple scales and classifiers along the architecture axis in Fig. 6b. The results indicate that the model rankings measured by the accuracy drop change for different 456 scales and shifts. For example, while the rankings remain consistent for the cartoon style (right) for 457 all scales, the model rankings change significantly for the painting style shift: Here, ViT outperforms 458 the other models on a lower scale but performs worse on large shift scales. Varying rankings also 459 occur for other shifts (see Fig. 9 in the supplementary). To validate the observation of changed model 460 rankings, we also evaluate multiple corruption levels of an examplary ImageNet-C corruption and 461 show the results in Fig. 23 in the supplementary. We conclude from this observation that the average 462 accuracy drop and the accuracy drops at specific nuisance scales do not always indicate the same 463 model behavior, which provides experimental evidence for the need for a multi-scale robustness 464 benchmarking dataset and adequate metrics.

- 465 Model failure points differ across different types of shifts. A failure point captures at which 466 scale a model fails for the first time. Comparing the failure point distribution of various models 467 largely differs for different shift types, as shown in Fig. 6c. We provide more results in Fig. 10 in 468 the supplementary. Weather shifts, such as snow, typically correspond to slight appearance changes 469 and mainly add a disturbance factor or occlusions to the image. Therefore, the failure rate increases 470 gradually compared to some style shifts, for which models tend to fail more abruptly at a specific scale, as, e.g., for the cartoon style at scale s = 1.5. An exemplary explanation for the abrupt shift 471 for the cartoon shift might be the wrong classification of a class as the ImageNet class *comic book*. 472
- 473 The relation between ID and OOD accuracy depends on the considered nuisance factor and 474 its scale. Miller et al. (2021) formalize the positive correlation between ID and OOD accuracy-475 classifiers tend to have a better OOD accuracy if they perform better on the training data ("Accuracy-476 on-the-line" phenomenon). To analyze the linear relation between ID and OOD accuracy for our 477 benchmark, we compute the slope of the linear fit between ID and OOD accuracies of 16 ImageNettrained models. Miller et al. (2021) have already shown that the slope varies for different datasets. In 478 Fig. 7, we further observe that not only the considered shift but also its severity influence the slope 479 of the linear fit. Refer to Appendix A.2.3 for the test statistics. We believe using our benchmark to 480 investigate this relation more extensively is an interesting direction for future work. 481
- Transformers with modern training recipes outperform modern CNNs across all shift sever ities. We present the average accuracy drops of various models with the same training data and a
   comparable number of parameters in Fig. 6a (*left*). DeiT3 consistently achieves the highest robust ness on our benchmark, increasing the gap towards DeiT and ViT for stronger shifts. Interestingly,
   ResNet-152 is more robust than the standard ViT variant, but ConvNeXt outperforms the ResNet-

<sup>486</sup> 152 architecture. A modern CNN (ConvNext) outperforms vision transformers (ViT,DeiT) of the <sup>487</sup> same size but it is less robust than a transformer with modern training recipes (DeiT3), despite <sup>488</sup> having a higher ID accuracy. This observation is in line with the performance on ImageNet-R. How-<sup>489</sup> ever, our benchmark shows that the gap between ConvNext and DeiT3 does not increase for stronger <sup>490</sup> shifts. We can observe that this behavior is not consistent for all shifts. Consider, *e.g.*, the failure <sup>491</sup> point distribution in Fig. 6c (*Painting Style*), where DeiT3 has a gradually increasing failure point <sup>492</sup> rate, while ConvNext depicts a sharp increase for scale s = 1.5.

Self-supervised pre-training improves the OOD robustness. To study the impact of the pre-training paradigm, we compare different learning objectives with the same ViT-B backbone and the same training data in Fig. 6a (*right*). We consider both the supervised and self-supervised (MAE, DINOv1, and MoCov3) paradigms. Using a self-supervised objective for pre-training followed by a fine-tuning protocol results in a better robustness for the same training data and model size. Considering the rCE metric in Table 1, the fine-tuned DINOv1 model achieves the best performance.

499 Diffusion classifiers are less robust than discriminative models. In addition, we also compare the 500 robustness of an ImageNet-trained diffusion classifier (Li et al., 2023b) on our benchmark. Due to 501 the heavy computational cost, we evaluate the accuracy drop of the DiT-based diffusion classifier for 1k images on a subset of our dataset (around 12k images) for the snow and the cartoon style shift. 502 We apply the L1 loss computation strategy as proposed by Li et al. (2023b) since it results in the best 503 performance. We compute the average accuracy drops as 0.106 / 0.07 / 0.05 for DiT / supervised 504 ViT / MAE. Comparing on the smaller dataset with discriminative models, the diffusion classifier 505 demonstrates a lower robustness on the evaluated shifts than the compared discriminative models 506 despite having substantially more parameters. The gap is increasing for larger severity levels. 507

More training data improves the robustness. In Fig. 6a (*right*), we observe that more training data
benefits OOD robustness for all scales. For example, compared with the supervised model trained
on IN1k, pre-training on IN21k positively impacts the OOD robustness for all scales. This might be
explained by the fact that the tested distribution is less OOD for the model (Miller et al., 2021).

512 In summary, we show that benchmarking with generative continuous shifts allows systematically 513 studying the model robustness via easily scalable synthetic data. Our study underscores that con-514 sidering multiple-scale nuisance shifts provides a more nuanced view of the model robustness, as 515 the performance drops can vary across different nuisance shifts and scales. Besides, the relation between ID and OOD accuracy not only depends on the considered nuisance factor but also on its 516 severity. Therefore, instead of aggregating the robustness evaluation into a single metric, we mo-517 tivate the community to report the accuracy with different shift scales and the failure points for a 518 more comprehensive understanding of model robustness. 519

520 521

522

#### 5 CONCLUSION

523 The key advantage of using generative models for benchmarking is the ability to perform diverse nuisance shifts in a controlled and scalable way. This work filled a gap in generative benchmarking 524 by introducing CNS-Bench, an evaluation method that performs diverse, realistic, fine-grained, and 525 continuous nuisance shifts at multiple scales. We further added a new dimension for benchmarking 526 robustness by introducing the concept of failure points. Our systematic evaluation of classifiers 527 revealed new insights along three axes (architecture, number of parameters, pre-training paradigm 528 and data) and demonstrated the importance of continuous shifts in assessing the model robustness. 529 Furthermore, we studied the necessity of removing out-of-class samples when benchmarking with 530 diffusion-generated images. Limitations and Future Work. While our approach allows for diverse 531 continuous nuisance shifts, it does not eliminate all confounders, meaning failures cannot always 532 be solely attributed to the targeted nuisance concept. This highlights an inherent challenge for 533 generative benchmarking approaches, and future advances in generative models could help mitigate 534 these confounding factors. Additionally, while we have carefully addressed this issue in our work, we acknowledge that using generated images can lead to biases arising from the real vs. synthetic 535 distribution shift. 536

We hope this benchmark can encourage the community to continue working on more high-quality
 generative benchmarks and to adopt generated images as an additional source for systematically
 evaluating the robustness of vision models.

All steps of our benchmarking pipeline are reproducible: We provide our datasets and implementation as part of the supplementary material, which includes code to reproduce training of LoRA adapters, generation of images, filtering, and evaluation of all classifiers. We also include all evaluated classification results for all images of the dataset in the shared code. All classifiers are evaluated in a standardized way using the *easyrobust* (Mao et al., 2022) framework.

The supplementary material contains more details about the implementation, the computation of metrics, the labeling, and the filtering strategies.

- 550 We also refer to our datasheet in Appendix B.
- 551 552

553

554

555 556

558

559

560

561

562

563

568

569

570

571

572

573

577

578

579 580

581

582

583

584

585

586

588

589

590

#### References

- Anastasios N Angelopoulos, Stephen Bates, Clara Fannjiang, Michael I Jordan, and Teodor Zrnic. Prediction-powered inference. *Science*, 382:669–674, Nov 2023.
- 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. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019.
- Stefan Andreas Baumann, Felix Krause, Michael Neumayr, Nick Stracke, Vincent Tao Hu, and Björn Ommer. Continuous, subject-specific attribute control in t2i models by identifying semantic directions, 2024. URL http://arxiv.org/abs/2403.17064.
- Florian Bordes, Shashank Shekhar, Mark Ibrahim, Diane Bouchacourt, Pascal Vincent, and Ari Mor Florian Bordes, Shashank Shekhar, Mark Ibrahim, Diane Bouchacourt, Pascal Vincent, and Ari Mor Pug: Photorealistic and semantically controllable synthetic data for representation learning.
   Advances in Neural Information Processing Systems, 36, 2024.
  - Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021.
  - X. Chen, S. Xie, and K. He. An empirical study of training self-supervised vision transformers. In 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 9620–9629, Los Alamitos, CA, USA, oct 2021. IEEE Computer Society.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE conference on computer vision and pattern recognition*, pp. 248–255. IEEE, 2009.
  - Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE* Signal Processing Magazine, 29(6):141–142, 2012.
  - 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. In *ICLR*, 2020.
  - Nathan Drenkow, Numair Sani, Ilya Shpitser, and Mathias Unberath. Robustness in deep learning for computer vision: Mind the gap? *CoRR*, abs/2112.00639, 2021. URL https://arxiv.org/abs/2112.00639.
  - A. Dutta, A. Gupta, and A. Zissermann. VGG image annotator (VIA). http://www.robots.ox.ac.uk/ vgg/software/via/, 2016. Version: X.Y.Z, Accessed: 2024-05-12.
- Abhishek Dutta and Andrew Zisserman. The VIA annotation software for images, audio and video.
   In Proceedings of the 27th ACM International Conference on Multimedia, MM '19, New York,
   NY, USA, 2019. ACM. ISBN 978-1-4503-6889-6/19/10. doi: 10.1145/3343031.3350535. URL
   https://doi.org/10.1145/3343031.3350535.

594 595	Bradley Efron. Tweedie's formula and selection bias. <i>Journal of the American Statistical Associa-</i> <i>tion</i> , 106(496):1602–1614, 2011.
597 598 599	Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. <i>International journal of computer vision</i> , 88: 303–338, 2010.
600 601	Lijie Fan, Kaifeng Chen, Dilip Krishnan, Dina Katabi, Phillip Isola, and Yonglong Tian. The scaling law of synthetic images for model training, for now. In <i>CVPR</i> , 2024.
602 603 604 605	Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In <i>The Eleventh International Conference on Learning Representations</i> , 2023.
606 607 608	Rohit Gandikota, Joanna Materzynska, Tingrui Zhou, Antonio Torralba, and David Bau. Concept sliders: LoRA adaptors for precise control in diffusion models, 2023. URL http://arxiv.org/abs/2311.12092.
609 610 611	Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III au2, and Kate Crawford. Datasheets for datasets, 2021.
612 613 614	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, 2022. URL http://arxiv.org/abs/1811.12231.
615 616 617	Trevor Hastie, Robert Tibshirani, and Jerome Friedman. <i>The Elements of Statistical Learning: Data Mining, Inference, and Prediction.</i> Springer Series in Statistics. Springer New York, NY, 2 edition, 2009.
619 620	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog- nition. In <i>CVPR</i> , 2016.
621 622	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In <i>CVPR</i> , 2022a.
624 625 626	Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiao- juan Qi. Is synthetic data from generative models ready for image recognition? <i>arXiv preprint</i> <i>arXiv:2210.07574</i> , 2022b.
627 628	Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common cor- ruptions and perturbations. In <i>International Conference on Learning Representations</i> , 2018.
629 630 631 632	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 <i>ICCV</i> , 2021a.
633 634	Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In <i>CVPR</i> , 2021b.
635 636 637 638	Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 33, pp. 6840–6851. Curran Associates, Inc., 2020.
639 640 641	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021.
642 643 644 645 646	Badr Youbi Idrissi, Diane Bouchacourt, Randall Balestriero, Ivan Evtimov, Caner Hazirbas, Nico- las Ballas, Pascal Vincent, Michal Drozdzal, David Lopez-Paz, and Mark Ibrahim. Imagenet- x: Understanding model mistakes with factor of variation annotations. <i>arXiv preprint</i> <i>arXiv:2211.01866</i> , 2022.

647 Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *CVPR*, pp. 18963–18974, 2022.

648 Alexander C. Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your diffu-649 sion model is secretly a zero-shot classifier, 2023a. 650 Alexander C Li, Mihir Prabhudesai, Shivam Duggal, Ellis Brown, and Deepak Pathak. Your dif-651 fusion model is secretly a zero-shot classifier. In Proceedings of the IEEE/CVF International 652 Conference on Computer Vision, pp. 2206–2217, 2023b. 653 654 Xiaodan Li, Yuefeng Chen, Yao Zhu, Shuhui Wang, Rong Zhang, and Hui Xue. Imagenet-e: Bench-655 marking neural network robustness via attribute editing. In CVPR, pp. 20371–20381, 2023c. 656 657 Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Doll'a r, and C. Lawrence Zitnick. Microsoft COCO: 658 common objects in context. CoRR, abs/1405.0312, 2014. URL http://arxiv.org/abs/ 659 1405.0312. 660 661 Jiang Liu, Chen Wei, Yuxiang Guo, Heng Yu, Alan Yuille, Soheil Feizi, Chun Pong Lau, and 662 Rama Chellappa. Instruct2attack: Language-guided semantic adversarial attacks. arXiv preprint 663 arXiv:2311.15551, 2023. 664 665 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In CVPR, 2022. 666 667 Xiaofeng Mao, Yuefeng Chen, Xiaodan Li, Gege Qi, Ranjie Duan, Rong Zhang, and Hui Xue. 668 Easyrobust: A comprehensive and easy-to-use toolkit for robust computer vision. https:// 669 github.com/alibaba/easyrobust, 2022. 670 671 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 ICLR, 2021. 672 673 Jan Hendrik Metzen, Robin Hutmacher, N Grace Hua, Valentyn Boreiko, and Dan Zhang. Identifica-674 tion of systematic errors of image classifiers on rare subgroups. In Proceedings of the IEEE/CVF 675 International Conference on Computer Vision, pp. 5064–5073, 2023. 676 677 John Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar, 678 Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation between out-of-distribution and in-distribution generalization. In International Conference on 679 Machine Learning, 2021. 680 681 Mohammadreza Mofayezi and Yasamin Medghalchi. Benchmarking robustness to text-guided cor-682 ruptions. In 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops 683 (CVPRW), pp. 779–786, 2023. doi: 10.1109/CVPRW59228.2023.00085. 684 685 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for 686 editing real images using guided diffusion models. In CVPR, pp. 6038-6047, 2023. 687 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khalidov, 688 Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learn-689 ing robust visual features without supervision. TMLR, 2023. 690 691 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 692 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 693 models from natural language supervision. In ICML, 2021. 694 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers 695 generalize to imagenet? In ICML, 2019. 696 697 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Highresolution image synthesis with latent diffusion models. In CVPR, pp. 10684–10695, 2022. 699 Michelle Shu, Chenxi Liu, Weichao Qiu, and Alan Yuille. Identifying model weakness with adver-700 sarial examiner. In Proceedings of the AAAI conference on artificial intelligence, volume 34, pp. 11998-12006, 2020.

702 703 704	Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper- vised learning using nonequilibrium thermodynamics. In <i>Proceedings of the 32nd International</i> <i>Conference on International Conference on Machine Learning</i> , pp. 2256–2265. JMLR, 2015.
705 706 707	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In Interna- tional Conference on Learning Representations, 2021.
708 709 710	Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In <i>ICML</i> , 2021.
711 712	Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. In ECCV, 2022.
713 714 715	Joshua Vendrow, Saachi Jain, Logan Engstrom, and Aleksander Madry. Dataset interfaces: Diagnos- ing model failures using controllable counterfactual generation. <i>arXiv preprint arXiv:2302.07865</i> , 2023.
716 717 718 719 720	Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Ra- sul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/ huggingface/diffusers, 2022.
721 722	Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representa- tions by penalizing local predictive power. In <i>NeurIPS</i> , 2019.
723 724 725 726	Chenshuang Zhang, Fei Pan, Junmo Kim, In So Kweon, and Chengzhi Mao. Imagenet-d: Bench- marking neural network robustness on diffusion synthetic object. In <i>CVPR</i> , pp. 21752–21762, 2024.
727 728 729	Bingchen Zhao, Shaozuo Yu, Wufei Ma, Mingxin Yu, Shenxiao Mei, Angtian Wang, Ju He, Alan Yuille, and Adam Kortylewski. Ood-cv: A benchmark for robustness to out-of-distribution shifts of individual nuisances in natural images. In <i>ECCV</i> , 2022.
730 731 732 733	Bingchen Zhao, Jiahao Wang, Wufei Ma, Artur Jesslen, Siwei Yang, Shaozuo Yu, Oliver Zendel, Christian Theobalt, Alan Yuille, and Adam Kortylewski. Ood-cv-v2: An extended benchmark for robustness to out-of-distribution shifts of individual nuisances in natural images. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2024.
735 736 737 738 739 740 741 742 743	Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor- corrector framework for fast sampling of diffusion models. <i>NeurIPS</i> , 2023.
744 745 746 747	
748 749 750	
751 752 753 754	
755	

## 756 A APPENDIX

This appendix provides supplementary information that is not elaborated in our main paper: We will discuss more details about the benchmarking dataset, the filtering, and image generation strategy. Additionally, we will provide more results. We refer to GoogleDrive <sup>1</sup> folder for the benchmarking dataset, the filtering dataset, and the code to reproduce training, generation, filtering, and benchmarking.

763 764

765

A.1 BENCHMARK DETAILS

This section provides more details about the benchmarking dataset. We first discuss the presented metrics. We then provide examples of the dataset and its distribution.

768 769 A.1.1 Access to Benchmarking Dataset

The dataset contains 192, 168 images in total, with 32, 028 images per scale. We share all images on
Google Drive in the folder *CNS\_dataset*. Additionally, we add the anonymized metadata, including
the annotations, as a JSON file. We will publish the data on our own servers upon submission.
Currently, data is accessible via Google Drive and can be downloaded by running the following
commands:

```
775  # Install gdown package
776  pip install gdown
777
778  # Download image files from Google Drive (22GB)
779  gdown https://drive.google.com/uc?id=1GYQb1dHu26mcklnMHtiySjZpigu-CmqN
780  gdown https://drive.google.com/uc?id=1GYQb1dHu26mcklnMHtiySjZpigu-CmqN
780  gdown https://drive.google.com/uc?id=1q4aS6oCdZx3jry6y92i3ZgoyNmHXZeN9
781 
782  # Unzip the downloaded file
```

783 unzip benchmark.zip

784

796 797

#### 785 A.1.2 LIST OF SHIFTS AND EXAMPLE IMAGES 786

The results are averaged over the following 14 shifts: cartoon style, plush toy style, pencil sketch style, painting style, design of sculpture, graffiti style, video game renditions style, style of a tattoo, heavy snow, heavy rain, heavy fog, heavy smog, heavy dust, and heavy sand-storm (see examples in Fig. 33 and Fig. 34). We train the sliders using the prompt template "A picture of a {class} in {shift}".

- 792 A.2 MORE RESULTS
- A.2.1 LABELED DATASET
  - Fig. 8 presents more classifier results on the labeled dataset.
- 798 A.2.2 LARGE BENCHMARK

We provide a table of accuracies and accuracy drops for all evaluated models and scales and the average accuracy and accuracy drop in Table 3. As discussed in the main paper, we also provide the accuracy drops for the ResNet family in Fig. 12. Similar to the observations in Table 3, larger models result in a lower accuracy drop in average. Fig. 9 provides a more nuanced view on the model performances accross various architectures on all shifts. We provide functionality to load the classification results for all images of the dataset in the shared code. All results are computed in a standardized way using the *easyrobust* (Mao et al., 2022) framework.

The accuracies for the diffusion classifier are depicted in Fig. 21. Similar to the discussion in the paper, the results showcase that the generative classifier is less robust than a supervised classifier.

<sup>808</sup> 809

<sup>&</sup>lt;sup>1</sup>https://drive.google.com/drive/folders/1twcuMLBSvy\_lIRYhssiwivBv9eKsA\_ ul?usp=sharing



Figure 8: Classification accuracy on the labeled dataset for snow and cartoon shifts. The accuracy drops on the labeled dataset showcase that various classifiers have varying sensitivities on different shifts.



Figure 9: Accuracy drops of various architectures for all shifts. We present the accuracy drops for all shifts in our benchmark. The performance gaps vary for different shifts and scales.

We use the DiT-based diffusion classifier trained on ImageNet-1k using the available framework (Li et al., 2023a) and the default hyper-parameters with a resolution of 256. Due to high computational costs, we compute the results for 100 classes, four scales, for the snow and cartoon style shift, and for at most 20 seeds per class, scale, and shift.

#### A.2.3 STATISTICS OF ACCURACY-ON-THE-LINE COMPUTATION

Fig. 17 provides the p-values of the linear regression corresponding to the presented results in Fig. 7.

## A.3 FINE-TUNING WITH SYNTHETIC DATA

We fine-tune a ResNet-50 classifier using our synthetic data. We compare the original ImageNet-trained model to a model fine-tuned using 50% synthetic data and 50% ImageNet training data. As shown in Table 5, the fine-tuned model leads to improved performance on the shifted real-world dataset, without a significant decline on the original ImageNet dataset.



Figure 10: Failure point distributions for all shifts. We present the failure point distributions for all shifts in our benchmark. The failure point distributions vary for different shifts, quantifying the different ways the shifts influence model performance.

#### A.4 ACCURACY DROPS ON IMAGENET-C

To provide more evidence that the model rankings change for different scales, we consider 7 levels of contrast as a deterministic example corruption from ImageNet-C, based on the implementation of Hendrycks & Dietterich (2018). We present the accuracy drops for all corruption levels in Fig. 23. A global averaged metric fails to capture such variations.

#### A.5 IMPLEMENTATION DETAILS

In this section, we provide more implementation details about the dataset generation process.

#### A.5.1 IMPLEMENTATION DETAILS FOR IMAGE GENERATION

We use the standard diffusers (von Platen et al., 2022) pipeline for Stable Diffusion 2.0, the DDIM
(Song et al., 2021) sampler with 100 steps and a guidance scale of 7.5, seeds ranging from 1 to 50.

A.5.2 Ablation of Image Generation

915 We ablate how the number of classes influences the robustness evaluations in Fig. 25. For a more 916 efficient computation, we use the UniPCMultistepScheduler sampler with 20 steps (Zhao 917 et al., 2023). In addition to 100 sliders for 14 shifts, we also publish the sliders for all 1000 ImageNet classes for the shifts snow and cartoon.



Figure 11: Failure point distributions for all shifts. We present the failure point distributions for all shifts in our benchmark along the model size axis of DeiT3. The failure point distributions vary for different shifts.



Figure 12: **Robustness evaluation for ResNet model family.** We compute the accuracy drops for all scales when varying the model size for a set of ResNet models. Larger models result in a better OOD robustness.

A.5.3 TEXT-BASED CONTINUOUS SHIFT

A naive approach for realizing continuous shifts involves computing the difference between two corresponding CLIP embeddings. We explored this strategy following the implementation of Baumann et al. (2024), but we did not achieve robust nuisance shifts for the variety of classes we considered and we present some examples in Fig. 26. We achieve reasonable results for some classes (*e.g.*, upper row). However, we observe the following issues arising from this strategy: (1) The semantic structures clearly change, which involves other factors of variation. This does not allow the computation of a failure point along one sliding trajectory. (2) depicted in middle row: For some classes,



Figure 13: Accuracy drops with confidence intervals. The accuracy drops are depicted for the three evaluation axes averaged over all shifts including the confidence interval of the accuracy computation.



Figure 14: Accuracy drops with confidence intervals. The accuracy drops are depicted for the three shifts along the model axes including the one-sigma confidence interval of the accuracy computation. The results show that some ranking changes are statistically stable.



Figure 15: Delta CLIP score with confidence. The  $\Delta$  CLIP score is plotted for various scales and averaged over all shifts including the standard deviation. The deviations are high. However, the can be attributed to the fact the range of the shift alignments varies for different shift types.



Figure 16: **Example of an incorrect CLIP alignment measure.** The CLIP alignment is applied as a measure to quantify the strength of the applied nuisance shift. However, this metric is not always correctly capturing the shift. Images with increasing slider scales and painting shift are represented. However, the alignment to the prompt "a picture in painting style" drops, as the  $\Delta$  CLIP difference to the first image depicted as the image titles demonstrates.

Table 2: mCE and mean rCE. We present the mean corruption error and the mean relative corrup-tion error for all evaluated models.

1028			
1029	Model	CE	rCE
1030	alexnet	1.000	1.000
1001	clip_resnet101	0.532	0.563
1031	clip_resnet50	0.715	0.587
1032	clip_vit_base_patch16_224	0.420	0.230
1002	clip_vit_base_patch32_224	0.487	0.591
1033	clip_vit_large_patch14_224	0.445	0.228
1004	clip_vit_large_patch14_336	0.419	0.274
1034	convnext_base.fb_in1k	0.359	0.686
1035	convnext_large.tb_in1k	0.354	0.672
	convnext_small.fb_in1k	0.353	0.609
1036	convnext_tiny.to_in1k	0.393	0.809
1037	convnextv2_base.tcinae_tt_in1k	0.322	0.553
1037	convnextv2_intge.icinae_it_in1k	0.283	0.555
1038	deit3 base patch16 224 fb in1k	0.396	0.508
1000	deit3 huge patch14 224 fb in1k	0.353	0.583
1039	deit3_large_patch16_224.fb_in1k	0.382	0.574
1040	deit3_medium_patch16_224.fb_in1k	0.387	0.758
1040	deit3_small_patch16_224.fb_in1k	0.400	0.747
1041	deit_base_patch16_224.fb_in1k	0.437	0.746
1040	dino_vit_base_patch16	0.504	0.851
1042	dinov1_vit_base_patch16	0.381	0.636
1043	dinov2_vit_base_patch14	0.350	0.524
	dinov2_vit_base_patch14_reg	0.311	0.456
1044	dinov2_vit_giant_patch14	0.321	0.431
10/15	dinov2_vit_giant_patch14_reg	0.311	0.426
1045	dinov2_vit_large_patch14	0.298	0.349
1046	dinov2_vit_large_patch14_reg	0.296	0.370
1017	dinov2_vit_small_patch14	0.351	0.639
1047	dinov2_vit_smail_patch14_reg	0.330	0.627
1048	mae_vit_base_patch14	0.380	0.752
10-10	mae_vit_large_patch16	0.305	0.571
1049	mac_vit_lage_paten10	0.320	0.669
1050	respect101.a1 in1k	0.491	0.842
1050	resnet152.a1_in1k	0.498	0.790
1051	resnet18.a1_in1k	0.493	0.954
	resnet34.a1_in1k	0.440	0.843
1052	resnet50.a1_in1k	0.485	0.945
1053	vit_base_patch16_224.augreg_in1k	0.569	0.926
1055	vit_base_patch16_224.augreg_in21k_ft_in1k	0.460	0.722
1054	vit_base_patch16_clip_224.openai_ft_in1k	0.282	0.482
1055			
1056			
1057	$10^{-1}$ 0.5 1.0 1.5 2.0 2.5		
1058			



Figure 17: p-values of the linear regressions corresponding to the plot in Fig. 7: The p-value is smaller than 0.05 for most scales and shifts, providing evidence for the statistical significance of our statements.

the naive approach is very unstable, resulting in OOD samples that do not represent realistic images. We did not reach significantly better results when applying a delayed sampling technique for the delta embedding. (3) depicted in the bottom row: Applying the delta in text-embedding space does not always result in a consistent increase of the considered shift. 

A.5.4 EVALUATION OF THE APPLIED SLIDER SHIFT

We evaluate whether our sliders always increase the shift, as measured by the  $\Delta$  CLIP score. For this purpose, we compute the  $\Delta$  CLIP score scores when increasing the slider scale by 0.5. Here,

1083	Shift Scale												
1005	Accuracy Accuracy Drop												
1084	model	0	0.5	1	1.5	2	2.5	avg	1	1.5	2	2.5	avg
1095	clip_resnet50	0.81	0.81	0.8	0.78	0.74	0.67	0.77	0.01	0.03	0.07	0.14	0.04
COUL	clip_resnet101	0.86	0.86	0.85	0.83	0.81	0.74	0.82	0.01	0.03	0.06	0.12	0.04
1086	clip_vit_base_patch16_224	0.87	0.88	0.88	0.87	0.86	0.81	0.86	-0.00	0.01	0.02	0.06	0.02
1007	clip_vit_base_patch14_224	0.87	0.87	0.80	0.85	0.85	0.77	0.84	-0.00	0.02	0.04	0.1	0.05
1087	clip vit large patch14 336	0.87	0.88	0.88	0.87	0.86	0.82	0.87	0.00	0.01	0.02	0.05	0.01
1088	convnext_tiny.fb_in1k	0.92	0.92	0.91	0.88	0.84	0.77	0.87	0.01	0.04	0.08	0.15	0.05
1000	convnext_small.fb_in1k	0.92	0.93	0.92	0.89	0.86	0.8	0.89	0.01	0.03	0.07	0.13	0.04
1089	convnext_base.fb_in1k	0.93	0.93	0.92	0.89	0.85	0.79	0.89	0.01	0.03	0.07	0.13	0.04
1090	convnext_large.fb_in1k	0.93	0.92	0.92	0.89	0.86	0.8	0.89	0.01	0.04	0.07	0.12	0.04
	convnextv2_base.fcmae_ft_in1k	0.93	0.93	0.92	0.9	0.87	0.82	0.9	0.01	0.04	0.07	0.12	0.04
1091	convnextv2_large.fcmae_ft_in1k	0.94	0.93	0.93	0.91	0.88	0.84	0.91	0.01	0.03	0.05	0.1	0.03
1092	convnextv2_huge.fcmae_ft_in1k	0.94	0.93	0.93	0.91	0.89	0.84	0.91	0.01	0.03	0.05	0.09	0.03
1002	doit3 base petch16 224 fb in1k	0.92	0.92	0.91	0.88	0.84	0.77	0.87	0.01	0.04	0.08	0.15	0.03
1093	deit3 medium patch16 224 fb in1k	0.91	0.91	0.9	0.88	0.84	0.79	0.87	0.01	0.03	0.07	0.12	0.04
100/	deit3_large_patch16_224.fb_in1k	0.91	0.91	0.9	0.88	0.85	0.8	0.88	0.01	0.03	0.06	0.12	0.04
1034	deit3_huge_patch14_224.fb_in1k	0.92	0.92	0.91	0.89	0.86	0.81	0.89	0.01	0.03	0.06	0.11	0.04
1095	deit_base_patch16_224.fb_in1k	0.9	0.9	0.89	0.87	0.83	0.76	0.86	0.01	0.04	0.08	0.15	0.05
1006	dino_lp_vit_base_patch16	0.9	0.9	0.89	0.85	0.8	0.71	0.84	0.01	0.05	0.1	0.19	0.06
1090	dinov1_ft_vit_base_patch16	0.91	0.91	0.90	0.88	0.84	0.84	0.87	0.01	0.03	0.07	0.04	asd
1097	dinov2_vit_small_patch14	0.92	0.92	0.91	0.89	0.86	0.81	0.89	0.01	0.03	0.06	0.11	0.04
1000	dinov2_vit_small_patch14_reg	0.93	0.93	0.92	0.9	0.87	0.81	0.89	0.01	0.03	0.06	0.11	0.04
1098	dinov2_vit_base_patch14	0.91	0.91	0.91	0.89	0.87	0.82	0.89	0.00	0.02	0.04	0.09	0.02
1099	dinov2_vit_base_patch14_leg	0.92	0.92	0.92	0.9	0.88	0.84	0.9	0.00	0.02	0.04	0.08	0.02
1100	dinov2_vit_large_patch14_reg	0.92	0.92	0.91	0.91	0.89	0.86	0.9	0.00	0.01	0.03	0.06	0.02
1100	dinov2_vit_giant_patch14	0.91	0.91	0.91	0.9	0.88	0.84	0.89	0.00	0.01	0.04	0.07	0.02
1101	dinov2_vit_giant_patch14_reg	0.92	0.92	0.91	0.9	0.88	0.85	0.9	0.00	0.01	0.03	0.07	0.02
	mae_vit_base_patch16	0.92	0.92	0.91	0.88	0.84	0.78	0.88	0.01	0.04	0.08	0.14	0.05
1102	mae_vit_huge_patch14	0.93	0.93	0.92	0.9	0.88	0.84	0.9	0.01	0.03	0.05	0.1	0.03
1103	mae_vit_large_patch16	0.93	0.92	0.92	0.9	0.87	0.83	0.9	0.01	0.03	0.05	0.1	0.03
1100	mocov3_vit_base_patch16	0.92	0.92	0.91	0.88	0.85	0.79	0.88	0.01	0.03	0.07	0.13	0.04
1104	respect24 a1 in1k	0.9	0.9	0.88	0.85	0.8	0.72	0.84	0.02	0.05	0.1	0.19	0.00
1105	respet50 a1 in1k	0.91	0.91	0.9	0.80	0.82	0.73	0.80	0.01	0.05	0.09	0.17	0.05
1105	resnet101.a1_in1k	0.9	0.9	0.88	0.85	0.8	0.73	0.84	0.02	0.05	0.1	0.17	0.06
1106	resnet152.a1_in1k	0.89	0.89	0.88	0.85	0.8	0.73	0.84	0.01	0.04	0.09	0.16	0.05
1107	vit_base_patch16_224.augreg_in1k	0.87	0.87	0.86	0.82	0.77	0.69	0.81	0.01	0.05	0.1	0.18	0.06
1107	vit_base_patch16_224.augreg_in21k_ft_in1k	0.9	0.9	0.89	0.86	0.82	0.75	0.85	0.01	0.04	0.08	0.15	0.05
1108	vit_base_patch16_clip_224.openai_ft_in1k	0.93	0.93	0.92	0.91	0.89	0.86	0.91	0.01	0.02	0.04	0.08	0.03
1100													
1109													
1110													
	0.8 0.5 1.0	1.5	2	.0	2.5								
1112													
4440										i I	-		
1113	ela ela												
1114	는 0.4 <b>- 1</b> - 1									1			
	ŭ l <b>u l</b>												

#### Table 3: Accuracy evaluations. We present the accuracies and accuracy drops of all evaluated classifiers.



supervised ImageNet-trained models for all evaluated shifts. The relation varies for different shifts and scales between 0.5 and 2.5.

0.2

0.0

the CLIP shift alignment increases for 73% of all cases for scales s > 0 and averaged over all shifts, demonstrating that increasing the slider weight results in a stronger severity of the desired shift.

#### A.5.5 IMPLEMENTATION DETAILS FOR BENCHMARKING

We provide the code for training the LoRA adapters and for performing the sliding. For benchmark-ing all vision models, we integrate our new benchmark and additional models in the easyrobust (Mao et al., 2022) framework. We provide all classification results for all images of the dataset together with the code and the data in the supplementary material.



Table 4: ImageNet validation accuracies and parameter count. One the left, we plot model
 accuracies on the ImageNet validation dataset for all evaluated classifiers. On the right, we present
 the parameter counts for the used architectures.

1191							
1192	Model	IN/val					
1193	clip_resnet101	58.00					
1194	clip_resnet50 clip_vit_base_patch16_224	55.00 67.70					
1105	clip_vit_base_patch32_224	62.60					
1195	clip_vit_large_patch14_224	75.00 76.30					
1196	convnext_base.fb_in1k	83.80					
1197	convnext_large.fb_in1k	84.30					
1198	convnext_small.tb_in1k convnext_tinv.fb_in1k	83.10 82.10		Model	Number of para	meters (in mil	lion)
1199	convnextv2_base.fcmae_ft_in1k	84.90		convnext_tiny			29
1000	convnextv2_huge.fcmae_ft_in1k	86.20 85.80		convnext_small			50 80
1200	deit3_base_patch16_224.fb_in1k	83.70		convnext_large			198
1201	deit3_huge_patch14_224.fb_in1k	85.10		convnextv2_base			89
1202	deit3_large_patch16_224.fb_in1k	84.60		convnextv2_huge			660 108
1202	deit3_small_patch16_224.fb_in1k	81.30		deit3_small			22
1203	deit_base_patch16_224.fb_in1k	81.80		deit3_medium			39
1204	dino_lp_vit_base_patch16	78.10		deit3_base			87
1205	dinov2-vit_base_patch16	82.49 84.50		deit3_huge			032 304
1206	dinov2_vit_base_patch14_reg	84.60		deit_base			87
1007	dinov2_vit_giant_patch14	86.60		vit_base			87
1207	dinov2_vit_large_patch14_leg	86.40		vit_large			307
1208	dinov2_vit_large_patch14_reg	86.70		resnet18			12
1209	dinov2_vit_small_patch14	81.40		resnet34			22
1210	mae_vit_base_patch16	83.70		resnet101			45
1210	mae_vit_huge_patch14	86.90		resnet152			60
1211	mae_vit_large_patch16	86.00					
1212	resnet101.a1.in1k	81.30					
1213	resnet152.a1_in1k	81.70					
1014	resnet18.a1_in1k	71.50					
1214	resnet50.a1_in1k	80.20					
1215	vit_base_patch16_224.augreg_in1k	76.80					
1216	vit_base_patch16_224.augreg_in21k_ft_in1k vit_base_patch16_clip_224.openai_ft_in1k	77.70 85.20					
1217	······································						
1218							
1219	0.9						
1220			ج (	).8			
1001	ğ 0.8 <b>*</b>	▲ ¥	urae		· · ·		
1221		·*  \	acc	).6 - SUP			
1222	6 0.7 - MoCo v3		OD	MoCo v.	3	k. T	
1223	MAE		0 (	).4 <b>† -</b> MAE			Ť
1224	0.6	×		→ DiffClas	s		*
1225	0.0 0.5 1.0 1.5	2.0 2.5		0.0 0.5	1.0 1.	.5 2.0	2.5
1226	scale				scale		
1227	(a) Accuracies for heavy snow	shift.		(b) Accuraci	es for cartoor	n style sh	ift.

Figure 21: **Comparison of DiT classifier.** We report the OOD accuracies for two shifts for the DiT classifier (Li et al., 2023b) and discriminative classifiers. All models were trained on ImageNet-1k and are evaluated on the same subset of our benchmark. The diffusion classifier performs worse than the discriminative models.

1233

1228

1234 One shift corresponds to a natural variation (snow), and the second shift corresponds to a style shift 1235 (cartoon style). (ii) We aim to find OOC samples that arise due to the application of the LoRA 1236 adapters. Therefore, we remove all images generated with a seed that results in a generated image 1237 that has a low CLIP text-alignment or is not classified classified correctly as the corresponding class even without the application of LoRA adapters. After removing such images, the labeling dataset 1238 consists of around 18k images. (iii) To reduce the labeling effort, we filter out all easy samples 1239 that (1) are correctly classified by DINOv2-ViT-L (Caron et al., 2021; Oquab et al., 2023) with a 1240 linear fine-tuned head and (2) one out of three classifiers (ResNet-50, DeiT-B/16, or ViT-B/16). (3) 1241 Additionally, the text alignment needs to be sufficiently high. (iv) Each hard image is labeled by



Table 5: **ImageNet-R performance after fine-tuning on our benchmark data**.ImageNet-R accuracy of the original ResNet-50 without fine-tuning and our model, fine-tuned on our benchmark.

Figure 22: Accuracy drops for three ImageNet-C corruptions and various architectures. The model rankings change for different corruptions, underlining the importance of the selection of the corruption types or nuisance shifts for benchmarking the OOD robustness. Additionally, it can also be observed that the accuracy drops at varying rates for different shifts.



Figure 23: Accuracy drops for contrast corruption. We report the accuracy drops for seven severities of the contrast corruption, as defined in (Hendrycks & Dietterich, 2018). The model rankings change for different scales.

1266 1267

1268

1270

1271 1272

1274

1276

1277

1278

1244

two human annotators. To increase the dataset quality, we include soft labels if the image partially
includes some characteristics of the class. So, each annotator can choose from the labels 'class',
'partial class properties', and 'not class'. An image is defined as an out-of-class sample if at least
one annotator considers the image as an OOC. For the remaining samples, an image is considered
IC (in-class) if at least one annotator labeled the image a clear sample of the corresponding class

For the pre-filtering strategy (ii) and for the selection of easy samples (iii), we compute text-alignment using CLIP score and we remove all samples that have a CLIP similarity  $s_{\text{CLIP-text-alignment}} > 24$ , which approximately includes 90% of all ImageNet validation images (Vendrow et al., 2023). We use the implementation in *torchmetrics* with VIT-B/16. After removing the easy samples in step (iii), 2.7k images remain for labeling. We use the VIA annotation tool (Dutta & Zisserman, 2019; Dutta et al., 2016) to create the annotations. Each image is labeled by two humans. In total, 14 graduate students are involved in the labeling process. For all participants, we ensure sufficient motivation and they receive detailed instructions on how to perform the labeling



Figure 24: ImageNet-R examples. Example images of one class where the shape and perspective significantly change.



(a) Accuracy over various scales.

(b) Failure point distribution (normalized over the sum of failure points).

Figure 25: Ablation of the number of ImageNet classes. We compare the accuracies and failure points averaged over the selected 100 classes and all 1000 ImageNet classes for two shifts (snow and cartoon style). We report the results with ResNet-50. The results indicate that the initial accuracy estimate is overestimated but the accuracy drops averaged over the two shifts are in line.

(the full set of instructions is provided in Fig. 32). We provide the filtering statistics in Table 6. An example screenshot of the labeling tool is visualized in Fig. 27. 

#### A.6.2 LABELING DATASET

We provide the images for labeling in the provided URL as well. There, we include all images and metadata that allow inferring the class of each image and the tag, whether it is labeled automatically or by a human. The statistics of the labeling dataset are shown in Fig. 28. 

A.7 USER STUDY 

We perform a user study on the final dataset using the same tooling as for the human labeling discussed in Appendix A.6 (iv). The user study includes 300 randomly sampled images and it is checked by two different individuals. In total, the user study involved seven people with different professions. 3 samples of our benchmark were considered as out-of-class samples, resulting in a ratio of 1% of failure cases with a margin of error of 0.5% for a one-sigma confidence level.

- A.8 APPLICATIONS OF TRAINED SLIDERS
- We can combing various sliders by simply adding the corresponding LoRA adapters. We show an example application in Fig. 29.



Figure 26: **Examples for text-based continuous shift.** The gradual increase can be successful. However, we observe that it fails for some classes (middle row) and is not consistently increasing (bottom row).

Table 6: Statistics of filtering process. We report the number of samples after various filtering stages. The stages are numbered according to the description in the main paper.

Scale	Stage (i)	Stage (ii)	Stage (iii)	Stage (iv)
0	4000	2966	2966	2966
0.5	4000	2966	2929	2955
1	4000	2966	2813	2906
1.5	4000	2966	2479	2740
2	4000	2966	2143	2498
2.5	4000	2966	1729	2110

1367

1368

# 1381 A.9 OOD-CV DETAILS

The Out-of-Distribution Benchmark for Robustness (OOD-CV) dataset includes real-world OOD
 examples of 10 object categories varying in terms of 5 nuisance factors: *pose, shape, context, texture,* and *weather.*

Generation of images for synthetic OOD-CV We generate the images for the synthetic OOD-CV dataset using a larger number of noise steps (85%) and more scale (between 0 and 3) since the classes occur more often in the dataset for training CLIP and Stable Diffusion. We use SD2.0 and not the dataset interfaces provided by Vendrow et al. (2023) since the class differences are less subtle and the samples of OOD-CV originate from two different datasets.

1392

1386

Training subset The OOD-CV benchmark provides a training subset of 8627 images. We train 1393 different state-of-the-art classifiers (i.e., ResNet-50 (He et al., 2016), ViT-B/16 (Dosovitskiy et al., 1394 2020), and DINO-v2-ViT (Oquab et al., 2023)) for classification. We finetune each baseline during 1395 50 epochs with an early stopping set to 5 epochs. In order to make baselines more robust, we apply 1396 standard data augmentation such as scale, rotation, and flipping during training. The training subset is composed of images originating from different datasets, notably ImageNet (Deng et al., 2009) 1398 and Pascal-VOC (Everingham et al., 2010). It is important to notice that the distribution of these 1399 two subsets is slightly different, with a higher data quality for the ImageNet subset and a lower quality for the latter subset (more noise, smaller objects, different image sizes). We visualize a few 1400 examples of the training data in Fig. 31. 1401

1402

**Test subset annotations** In the test subset provided in the benchmark dataset, only the coarse individual nuisance factors (*e.g.*, *weather*, *texture*) are provided. In our setup, we





(a) For the human labeling dataset.



Figure 28: **Statistics of labeling dataset.** We report the number of in-class, partially in-class, and out-of-class samples.

are interested in studying more fine-grained nuisance shifts, notably rain, snow, or fog. Hence, we had to assign some fine-grained annotation to all images containing *weather* nui-Hence, we assign a fine-grained annotation by computing the CLIP similar-sance shifts. ity to the following texts: "a picture of a {class} in {shift}", where class is the ground truth class and shift the nuisance shift candidate rain, snow, or fog and "a picture of a {class} without snow nor fog nor rain". By applying a softmax on the similarity scores with the previous texts, we can assign the fine-grained nuisance shift rain, snow, fog or unknown for each image. We show more statistics in Table 7. By checking the results visually, we observe that all fine-grained nuisance shifts align with human perception and have a tendency towards classifying samples as *unknown* as soon as there is a small doubt. Note that by applying the same strategies to our generated data, we obtain an accuracy close to 100%.

Nearest neighbor images of OOD-CV and CNS-Bench. To illustrate the realism of our gener ated image, we compute the nearest neighbours using cosine similarity with CLIP image embedding
 and we plot it in Fig. 30.

1458				
1459				
1460	1000			
1461			and the	
1462		Contraction of the second		
1463	A			
1464				
1404				
1400				
1400	and the second se		STA	
1407				and the second sec
1400				
1469				
1470				
14/1	100 C			
14/2				
1473			A P	
1474				
1475				
1476				
14//	Figure 29: Combination of Sli	iders. We ex	cemplarily s	show that sliders can be combined. Here, a
1478	snow slider (vertical axis) and a	cartoon slid	er (horizont	al axis) are linearly added for three scales.
1479				
1480				
1481				
1482				
1483	synthetic top1.1NN OOD-CV sample	synthetic top2	-NN synthe	elic top3-NN synthetic top5-NN
1484		- 7 -	d -	7 - 7 - 7
1485		DI	n at	
1486				
1487		State of the states		
1400				
1409	synthetic top1-NN	synthetic top2	-NN synthe	etic top3-4N synthetic top4-4N synthetic top5-NN
1490	000-CV sample		3	
1491				
1492		- Part all		
1493				
1494				
1495	Figure 30: <b>Closest synthetic sa</b>	mples to two	) example (	<b>DOD-CV images.</b> We find the top-5 nearest
1490	neighbours using cosine similar	ity with CLI	P image em	bedding.
1497	5 5	2	U	C
1490				
1499				
1500				
1501	Table 7: OOD-CV Statistics.	We report	the number	of images and accuracies for the weather
1502	subset.	-		-
1503				
1504		Shift	#images	Accuracy
1505		Snow	273	70.3
1507		Fog	24	62.5
1502		Rain	74	66.2
1500		Unknown	129	66.7
1510		Total	500	68.4

28



Figure 31: OOD-CV example images. We infustrate a set of example images from the training and the testing dataset of OOD-CV: (a-h) example from the training set, from ImageNet or Pascal-VOC.
(i-l) Some examples for weather nuisance shifts. In the training set, we observe that images from the Pascal-VOC subset are usually of lower quality (*e.g.*, cropping, occlusion, resolution) compared to the ImageNet subset. In the test set, we see that that not fully disentangled (*e.g.*, (j) is only partially visible, (k) is partially occluded).

1566		
1567		
1568		
1569	ا م ا	aling took for out of alage datastion
1570	Lap	lenny lask for out-or-class detection
1571	Motiva	tion: For benchmarking a classifier with synthetic images, we need to ensure that the generated images
1572	still cor	respond to the correct classes. To evaluate automatic filtering pipelines, we create a dataset with human
1573	labels.	The dataset includes generated images with various levels of snow or cartoon style.
1574	Task:	
1575	The go	al is to detect images that do not belong to the corresponding ImageNet class (given as title).
1576		
1577	Given	an image, your task is to select one of three labels:
1578	•	• You can clearly recognize the class
1579	•	partiv class:
1580		<ul> <li>Given the class label, the class seems to correspond to the image.</li> </ul>
1581		• You can recognize parts of the class but you are not very sure whether this is actually the class
1582		<ul> <li>You clearly see some characteristics of the class but it does not include all the important features</li> </ul>
1583	•	not class:
1584		<ul> <li>The considered image is clearly not the considered class.</li> </ul>
1585		
1586	The go	al is to check whether the objects in the image correspond to a class or not. The goal is not to check
1587	wnethe	er the samples look realistic.
1588	Every	class starts with one realistic example image, taken from ImageNet. This image needs to be labeled as
1589	well. S	ince the example is just one illustrative example, not depicting the diversity of the class, it is
1590	recom	mended to use Google picture search to get an intuition of how the object looks in case one is not familiar
1591	with the	e class. In the consecutive class complex will be similar. They are concreted with the same seed but with verying
1592	snow o	or cartoon levels.
1593	0.1011 0	bucket 📰 red fox
1594	Some	examples for class, partly class, and not class:
1595		
1590	1)	clase: This animal can be clearly described as a fox at first
1502	''	glance. Also, the bucket can be easily recognized.
1590		
1600		
1601		🗐 ladybug Granny Smith
1602		
1603	2)	partly class: The shape and size seems to fit a ladybug.
1604		However, the black dots are missing. The other picture might be
1605		a cartoon-like illustration of apples. However, this can be argued.
1606		
1607		
1608		- Sax harmerhead
1609		
1610	3)	not class: First example: This is supposed to be a say but it is
1611	5)	clearly not recognizable as a sax. Second example: There is not
1612		a single characteristic that resembles a hammerhead. It is very
1613		clearly not the class.
1614		
1615	Figure	32: Sat of instructions for labeling. Instructions provided to the human appointers to per-
1616	form th	be labeling of the out-of-class filtering dataset
1617	101111 11	to according of the out of endos intering unuser.
1618		



Figure 33: Example sliding for various nuisance shifts. We visualize six generated images wit the corresponding scales as 0, 0.5, 1, 1.5, 2, and 2.5.





## 1728 B DATASHEET

1730

1731 In the following, we answer the questions as proposed in Gebru et al. (2021).

1732

#### 1733 1734 B.1 MOTIVATION

For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

The dataset was created to evaluate the robustness of state-of-the-art models to specific continuous nuisance shifts. Current approaches are not scalable and often include only a small variety of nuisance shifts, which are not always relevant in the real world. More importantly, current benchmark datasets define binary nuisance shifts by considering the existence or absence of that shift, which may contradict their continuous realization in real-world scenarios.

Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

Until the acceptance of the paper, the specific details about the research group, their affiliations, and the entities they represent will remain anonymous.

1748 Who funded the creation of the dataset? If there is an associated grant, please provide the name of the 1749 grantor and the grant name and number.

- <sup>1750</sup> Until the acceptance of the paper, the specific details about funding will remain anonymous.
- 1751 1752
- 1753 B.2 COMPOSITION
- 1754

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?

The dataset consists of synthetic images that were generated using Stable Diffusion.

1759 How many instances are there in total (of each type, if appropriate)?

The dataset contains 192, 168 images in total, with 32, 028 for each of the six scales with 14 shifts.
Each shift has at least 5,000 images and 100 classes.

Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances because instances were withheld or unavailable).

The dataset contains the subset of images that were filtered using the selected filtering strategy. Originally, 420, 000 images were generated.

- What data does each instance consist of? "Raw" data (e.g., unprocessed text or images) or
  features? In either case, please provide a description.
- 1772 "Raw" synthetically generated data as described in the paper.
- 1774 **Is there a label or target associated with each instance?** If so, please provide a description.
- <sup>1775</sup> Yes, each image belongs to an ImageNet class and has a shift scale assigned to it.
- **Is any information missing from individual instances?** If so, please provide a description, explaining
- why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.
- No, for each instance, we give the class label, the scale of the shift, and the parameters used for
  generating this image. However, the class label might be erroneous in rare cases where the generated image corresponds to an out-of-class sample.

1782 Are relationships between individual instances made explicit (e.g., users with their tweets, 1783 songs with their lyrics, nodes with edges)? If so, please describe how these relationships are made 1784 explicit. 1785 Yes, the relationships in terms of class, random seed for generation, shift, and scale of shift are 1786 provided in the dataset. 1787 Are there recommended data splits (e.g., training, development/validation, testing)? If so, please 1788 provide a description of these splits, explaining the rationale behind them. 1789 1790 We offer a benchmark dataset specifically intended for testing the robustness of classifiers. There-1791 fore, we recommend utilizing the entire dataset provided as the test dataset. 1792 Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a 1793 description. 1794 We provided a dataset of generated images. While we apply a filtering strategy to reduce the number 1795 of out-of-class and unrealistic samples, we cannot guarantee that all images of the dataset represent a 1796 realistic and visually appealing realization of the considered class. We provide a statistical estimate 1797 of the number of failure samples in the paper. The data might also include the redundancies that underlie the image generation process of Stable Diffusion. 1799 Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., 1801 websites, tweets, other datasets)? If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset 1802 (i.e., including the external resources as they existed at the time the dataset was created); c) are there any 1803 restrictions (e.g., licenses, fees) associated with the use of these external resources? 1804 1805 The dataset is fully self-contained. 1806 Does the dataset contain data that might be considered confidential (e.g., data that is pro-1807 tected by legal privilege or by doctor-patient confidentiality, data that includes the content of 1808 individuals' non-public communications)? If so, please provide a description. 1809 No. 1810 1811 Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, 1812 or might otherwise cause anxiety? If so, please describe why. 1813 There is a small chance that our synthetically generated data can generate offensive images. How-1814 ever, we did not encounter any such sample during our extensive manual annotations. 1815 Does the dataset relate to people? If not, you may skip the remaining questions in this section. 1816 1817 1818 No. 1819 Does the dataset identify any subpopulations (e.g., by age, gender)? If so, please describe how these 1820 subpopulations are identified and provide a description of their respective distributions within the dataset. 1821 1822 N/A. Is it possible to identify individuals (i.e., one or more natural persons), either directly or indi-1824 rectly (i.e., in combination with other data) from the dataset? If so, please describe how. 1825 N/A. 1826 1827 Does the dataset contain data on individuals' protected characteristics (e.g., age, gender, race, 1828 **religion**, **sexual orientation**)? If so, please describe this data and how it was obtained. 1829 N/A. 1830 Does the dataset contain data on individuals' criminal history or other behaviors that would 1831 typically be considered sensitive or confidential? If so, please describe this data and how it was ob-1832 tained. 1833 1834 N/A. 1835

1836 1837	B.3 COLLECTION PROCESS
1838 1839 1840 1841	How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses)?
1842	N/A.
1843 1844 1845	What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)? How were these mechanisms or procedures validated?
1847 1848	We used Stable Diffusion 2.0 to generate all images. Images were generated using NVIDIA A100 and A40 GPUs.
1849 1850	If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?
1851 1852 1853	The dataset was filtered using a combinatorial selection approach using the alignment scores of DINOv2 and CLIP to the considered class.
1854 1855	Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?
1856 1857	The authors of the paper and other PhD students of the institute. They were not additionally paid for the dataset collection process.
1858 1859 1860	<b>Over what timeframe was the data collected? Does this timeframe match the creation time- frame of the data associated with the instances (e.g., recent crawl of old news articles)?</b> If not, please describe the timeframe in which the data associated with the instances was created.
1862	The images were generated and processed over a timeframe of four weeks.
1863 1864 1865	Were any ethical review processes conducted (e.g., by an institutional review board)? If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.
1866 1867 1868 1869 1870 1871	No ethical concerns. B.4 PREPROCESSING/CLEANING/LABELING
1872 1873 1874 1875 1876	Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remaining questions in this section.
1877 1878 1879 1880	Yes, cleaning of the generated data was conducted. The generated images underwent filtering to reduce the number of out-of-class samples using the proposed filtering mechanisms. Instances that did not meet these criteria were removed from the dataset. For a detailed description of the filtering process, please refer to the corresponding section in the paper.
1881 1882	Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the "raw" data.
1883 1884	The generated images remain in their original, unprocessed state and can be considered as "raw" data. However, we have not provided all the images that were filtered out.
1886 1887	Is the software used to preprocess/clean/label the instances available? If so, please provide a link or other access point.
1888 1889	Generating the images was performed using commonly available Python libraries. For annotating a subset of the dataset for filtering purposes, we have used the VIA annotation tool (Dutta & Zisserman, 2019; Dutta et al., 2016).

1890 B.5 Uses

**Has the dataset been used for any tasks already?** If so, please provide a description.

In our work, we demonstrate how this approach yields valuable insights into the robustness of stateof-the-art models, particularly in the context of classification tasks.

Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point.

- 1898 1899 We will provide a link that includes all relevant papers or systems.
- 1900 What (other) tasks could the dataset be used for?

1901 Our work showcases the capability of our dataset to enhance control over data generation, which 1902 is particularly evident through continuous shifts. However, its applicability extends beyond this 1903 demonstration. The dataset can be effectively utilized in various generation tasks that necessitate 1904 continuous parameter control. While we showcased its efficacy in providing insights for models 1905 tackling classification tasks, it can seamlessly extend to evaluate the robustness of state-of-the-art methods across diverse tasks such as segmentation, domain adaptation, and many others. This is possible by combining our approach with other modes of conditioning Stable Diffusion. In addi-1907 tion, our data can also be used for fine-tuning, which we also demonstrated in the supplementary 1908 material. 1909

Is there anything about the composition of the dataset or the way it was collected and cleaned
that might impact future uses? For example, is there anything that might cause the dataset to
be used inappropriately or misinterpreted (e.g., accidentally incorporating biases, reinforcing
stereotypes)?

Our dataset was synthesized using a generative model. It, therefore, likely inherits any biases for its generator. Similarly, filtering is performed by pre-trained models, which can indirectly also contribute to biases.

Are there tasks for which the dataset should not be used? If so, please provide a description.

No, there are no tasks for which the dataset should not be used. Our dataset aims to enhance model
robustness and provide deeper insights during model evaluation. Therefore, we see no reason to
restrict its usage.

1922 1923

1925

1924 B.6 DISTRIBUTION

Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created? If so, please provide a description.

1928 Yes, the dataset will be publicly available on the internet.

How will the dataset be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?

- 1932 In the future, we will distribute the dataset as a tarball on our servers.
- 1933 1934 When will the dataset be distributed?

1935 The dataset will be distributed upon acceptance of the manuscript. It is now available under the 1936 provided anonymized link.

Will the dataset be distributed under a copyright or other intellectual property (IP) license,
and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide
a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU.

1940 CC-BY-4.0.

1942 Have any third parties imposed IP-based or other restrictions on the data associated with the

**instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms.

1944
 1945
 1946
 No, there are no IP-based or other restrictions on the data associated with the instances imposed by third parties.

**Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.

- 1950 We are not aware of any export controls or other regulatory restrictions that apply to the dataset or 1951 to individual instances.
- 1952

1949

1953 B.7 MAINTENANCE

19541955 Who is supporting/hosting/maintaining the dataset?

The dataset is supported by the authors and their associated research groups. The dataset is hosted on our own servers.

How can the owner/curator/manager of the dataset be contacted (e.g., email address)?

The authors of this dataset will be reachable at their e-mail addresses: [undisclosed]. In addition, we will add a contact form, which will be made available on the website.

- **Is there an erratum?** If so, please provide a link or other access point.
- 1963 1964 If errors are found, an erratum will be added to the website.

Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-stances)? If so, please describe how often, when, and how updates will be provided.

Yes, updates will be communicated via the website. The dataset will be versioned.

If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a specific period of time and then deleted)? If so, please describe these limits and explain how they will be enforced.

- 1973 Our dataset does not relate to people.
- Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how.

No, older versions of the dataset will not be supported if the dataset is updated. We do not plan to extend or update the dataset. Any updates will be made solely to correct any hypothetical errors that may be discovered.

1980 If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for 1981 them to do so? If so, please provide a description. Will these contributions be made publicly available?

Yes, we provide all the necessary tools and explanations to enable users to build continuous shifts
for their own specific applications. Our dataset serves as a foundation to illustrate how it can be used
to evaluate current state-of-the-art methods. However, we are happy to centralize and showcase all
related work on our GitHub page that benefits from our method of generating data.

1986

1988

1987 B.8 AUTHOR STATEMENT OF RESPONSIBILITY

1989 The authors confirm all responsibility in case of violation of rights and confirm the license associated 1990 with the dataset and its images.

- 1991
- 1992
- 1993
- 1994
- 1995
- 1996
- 1997