

---

# Reproducibility Study of “Counterfactual Generative Networks”

---

Anonymous Author(s)

Affiliation

Address

email

## Reproducibility Summary

1

### 2 **Scope of Reproducibility**

3 In this work, we study the reproducibility of the paper *Counterfactual Generative Networks* (CGN) by Sauer and Geiger  
4 to verify their main claims, which state that (i) their proposed model can reliably generate high-quality counterfactual  
5 images by disentangling the shape, texture and background of the image into independent mechanisms, (ii) each  
6 independent mechanism has to be considered, and jointly optimizing all of them end-to-end is needed for high-quality  
7 images, and (iii) despite being synthetic, these counterfactual images can improve out-of-distribution performance of  
8 classifiers by making them invariant to spurious signals.

### 9 **Methodology**

10 The authors of the paper provide the implementation of CGN training in PyTorch. However, they did not provide code  
11 for all experiments. Consequently, we re-implemented the code for most experiments, and run each experiment on 1080  
12 Ti GPUs. Our reproducibility study comes at a total computational cost of 112 GPU hours.

### 13 **Results**

14 We find that the main claims of the paper of (i) generating high-quality counterfactuals, (ii) utilizing appropriate  
15 inductive biases, and (iii) using them to instil invariance in classifiers, do largely hold. However, we found certain  
16 experiments that were not directly reproducible due to either inconsistency between the paper and code, or incomplete  
17 specification of the necessary hyperparameters. Further, we were unable to reproduce a subset of experiments on a  
18 large-scale dataset due to resource constraints, for which we compensate by performing those on a smaller version of  
19 the same dataset with our results supporting the general performance trend.

### 20 **What was easy**

21 The original paper provides an extensive appendix with implementation details and hyperparameters. Beyond that, the  
22 original code implementation was publicly accessible and well structured. As such, getting started with the experiments  
23 proved to be quite straightforward. The implementation included configuration files, download scripts for the pretrained  
24 weights and datasets, and clear instructions on how to get started with the framework.

### 25 **What was difficult**

26 Some of the experiments required severe modifications to the provided code. Additionally, some details required for the  
27 implementation are not specified in the paper or inconsistent with the specifications in the code. Lastly, in evaluating  
28 out-of-distribution robustness, getting the baseline model to work and obtaining numbers similar to those reported in  
29 the respective papers was challenging, partly due to baseline model inconsistencies within the literature.

### 30 **Communication with original authors**

31 We have reached out to the original authors to get clarifications regarding the setup of some of the experiments, but  
32 unfortunately, we received a late response and only a subset of our questions was answered.

# 1 INTRODUCTION

Despite the considerable popularity of deep learning models within the field of artificial intelligence, recent literature suggests that these networks have a tendency of learning simple correlations that perform well on a benchmark dataset, instead of more complex relations that generalize better [1, 17, 21]. This phenomenon, which is referred to as shortcut learning by Geirhos et al. [10], makes these models more sensitive to input perturbation and unseen input contexts.

In order to enhance the robustness and interpretability of classifiers, Sauer and Geiger [22] introduce the idea of a *Counterfactual Generative Network* (CGN). Using appropriate inductive biases to disentangle separate components within the input images, such as object shape, object texture, and background, this model is capable of augmenting training data with generated counterfactual images. The authors claim that, using this model, they were able to improve out-of-distribution (OOD) robustness with only a marginal performance decrease for the original classification task.

In this work, we aim to reproduce their findings, verify their claims, and perform additional experimental results to provide further evidence to support their claims. In summary, this work makes the following contributions:

- We reproduce the main experiments conducted by Sauer and Geiger [22] to identify which parts of the experimental results supporting their claims can be reproduced, and at what cost in terms of resources (e.g., computational cost, development effort, and communication with the authors).
- We improve the performance consistency of the CGN during training.
- We extend upon the work of Sauer and Geiger by empirically analyzing the decisions made by classifiers based on their proposed model. Based on this analysis, we propose a method to quantify the robustness of such classifiers against spurious correlations.

## 1.1 Scope of Reproducibility

Distinguishing between spurious and causal correlation is an active topic in causality research [15, 18]. One central principle in causal inference is the assumption of independent mechanisms (IMs), which states that a causal generative process is composed of autonomous modules that do not influence each other [19, 22, 24]. The CGN introduced in the original paper exploits this idea to decompose the image generation process into three IMs, each controlling one factor of variation (FoV), namely the shape, texture, and background. Using this, the authors take a step towards more robust and interpretable classifiers that explicitly expose the causal structure of the classification task. In this reproducibility study, our main goal is to verify the following claims of the original paper:

- **High-Quality Counterfactuals (HQC):** By exploiting proper inductive biases, the CGN is able to reliably learn the independent mechanisms, which allow for the generation of high-quality counterfactual images by disentangling the shape, texture and background of the image.
- **Inductive Bias Requirements (IBR):** Each independent mechanism has to be considered, and jointly optimizing all of them end-to-end is needed for high-quality images.
- **Out-of-Distribution Robustness (ODR):** Despite being synthetic, the counterfactual images can improve out-of-distribution performance of classifiers by making them invariant to spurious signals.

The remainder of this work is structured as follows. In Section 2, we introduce the model proposed in the original paper to provide the reader with the required background knowledge. Section 3 then summarizes our approach to reproduce the original paper. Subsequently, Section 4 presents the replicated results and compares them to the original paper. Section 5 concludes this work by discussing our experience with reproducing the research by Sauer and Geiger [22].

# 2 COUNTERFACTUAL GENERATIVE NETWORK

The counterfactual generative network is a manifestation of a structural causal model (SCM) for the task of image classification [22]. It decomposes the image generation process into four IMs whose losses are jointly optimized in an end-to-end manner. An overview of the CGN architecture is shown in Appendix A.

**Shape mechanism:** The shape mechanism  $f_{shape}$  captures the shape as a binary mask  $m$ , where 1 corresponds to the object and 0 to the background. For this purpose, it first samples a pre-mask  $\tilde{m}$  with exaggerated features from a

77 fine-tuned BigGAN [4], and extracts the binary mask using a pretrained U2-Net [20]. The shape loss  $\mathcal{L}_{shape}$  comprises  
 78 (1) the *pixelwise binary entropy* of the mask, and (2) the mask loss:

$$\mathcal{L}_{mask}(\mathbf{m}) = \mathbb{E}_{p(\mathbf{u}, y)} \left[ \max \left( 0, \tau - \frac{1}{N} \sum_{i=1}^N m_i \right) + \max \left( 0, \frac{1}{N} \sum_{i=1}^N m_i - \tau \right) \right]. \quad (1)$$

79 The pixelwise binary entropy forces the output to be close to either 0 or 1, whereas the mask loss discourages trivial  
 80 solutions that are outside the interval defined by  $\tau$ .

81 **Texture mechanism:** The texture mechanism  $f_{text}$  generates the texture of the object. For MNIST, Sauer and Geiger  
 82 use an additional layer that divides its input into patches and randomly rearranges them. In contrast, for ImageNet, they  
 83 sample patches from the regions where the mask values are the highest and concatenate them into a patch grid  $pg$ . This  
 84 mechanism is optimized by minimizing the perceptual loss between the foreground  $\mathbf{f}$  and the patch grid  $pg$ . As such,  
 85 the background gradually transforms into object texture during training.

86 **Background mechanism:** The background mechanism  $f_{bg}$  models the background  $\mathbf{b}$  of the image. It removes the  
 87 object from the output of the BigGAN backbone and inpaints it using U2-Net by *minimizing* the predicted saliency.  
 88 Because there is no need for a globally coherent background in the MNIST setting, the MNIST variant of the CGN  
 89 includes a second texture mechanism rather than a dedicated background mechanism.

90 **Composer:** The composer  $C$  combines the output of the aforementioned mechanisms into a single composite image

$$\mathbf{x}_{gen} = C(\mathbf{m}, \mathbf{f}, \mathbf{b}) = \mathbf{m} \odot \mathbf{f} + (1 - \mathbf{m}) \odot \mathbf{b}, \quad (2)$$

91 where  $\mathbf{m}$  is the mask,  $\mathbf{f}$  is the foreground,  $\mathbf{b}$  is the background, and  $\odot$  is the Hadamard product. To optimize this  
 92 mechanism, Sauer and Geiger use an external conditional GAN (cGAN) that generates pseudo-ground-truth images  
 93  $\mathbf{x}_{gt}$  from the same noise  $\mathbf{u}$  and label  $y$  that is fed into the aforementioned mechanisms of the CGN. Using this, they  
 94 minimize the reconstruction loss  $\mathcal{L}_{rec}$  between the composite image  $\mathbf{x}_{gen}$  and the pseudo-ground-truth image  $\mathbf{x}_{gt}$ .

95 During training, each independent mechanism learns a class-conditional distribution over shapes, textures, or back-  
 96 grounds. It can then generate counterfactual images by randomizing the noise  $\mathbf{u}$  and label input  $y$  for each mechanism.  
 97 A more detailed explanation regarding the purpose of these counterfactual images and the connection with explainable  
 98 artificial intelligence (XAI) can be found in Appendix B.

99 In order to encode invariance to spurious correlations, Sauer and Geiger train classifiers on generated counterfactual  
 100 data that retain the label from the shape with randomized texture and backgrounds. For MNISTs, they use a standard  
 101 CNN feature extractor followed by a single classification head. For ImageNet on the other hand, they use a CNN  
 102 backbone with three classifier heads: shape, texture, and background; each invariant to all but one factor of variation.  
 103 The final prediction is obtained by averaging the individual head predictions.

### 104 3 METHODOLOGY

105 The original implementation of the CGN is publicly available [23], but most of the experiments conducted in the  
 106 original paper to support their claims are not. Consequently, we use the authors’ code for the implementation of the  
 107 CGN, and re-implement the experiments and relevant evaluation metrics based on the descriptions provided in the paper.  
 108 Furthermore, we both improve and extend upon the work of Sauer and Geiger by providing additional experiments and  
 109 results. Because a description of the GAN used in the original paper was not provided, we use a DCGAN [14].

#### 110 3.1 Datasets

111 The experiments conducted in the original paper involve two tasks, namely generating counterfactual examples and  
 112 training a classifier to be invariant to spurious correlations. We follow the paper and reproduce their evaluations on  
 113 multiple datasets for each task. For both tasks, we present the relevant datasets and their main purpose in Table 1.  
 114 Due to resource constraints, running all experiments on full ImageNet (IN-1k) is infeasible. As a compromise, we use  
 115 ImageNet-mini (IN-Mini) [7], a small-scale variant of ImageNet. Although this dataset contains fewer samples, we  
 116 found it to be sufficient to reproduce the main findings of the original paper and verify their claims. Moreover, this  
 117 dataset includes the same classes as IN-1k and hence does not induce any decrease in difficulty of the classification task.

Table 1: **Datasets overview.** The datasets used for empirical evaluations across two tasks.

Task	Datasets	Number of samples			Classes	Description	URL
		Train	Test	Total			
Generating counterfactual samples	C-MNIST [2]	50k	10k	60k	10	Foreground colour as a spurious correlation	Link
	DC-MNIST	50k	10k	60k	10	Fore/background colour as spurious correlations	NA <sup>1</sup>
	W-MNIST	50k	10k	60k	10	In-the-wild background with texture colour	NA <sup>1</sup>
	IN-1k [6]	1M	100k	1.2M	1000	Large-scale evaluation	Link
	IN-mini [7]	35k	4k	39k	1000	Small-scale evaluation	Link
Training invariant classifiers	MNISTs	50k	10k	60k	10	Test different granularities of invariance	Link
	Cue-conflict [8]	NA	1280	1280	16	Tests shape-texture disentanglement	Link
	IN-9 variants [29]	~45k	~4k	~50k	9	Tests background-invariance	Link

## 3.2 Hyperparameters

In order to match the original experiments as closely as possible, we used the same hyperparameters as the authors of the original paper whenever they were specified in the article. If the required hyperparameters for the experiments were not mentioned in the original paper, we relied on the default parameters given in the configuration files of the original implementation. In this case, we assume that these default parameters correspond to the parameters used for the described experiments.

## 3.3 Experimental setup and evaluation metrics

Our experimental setup is largely based on the description provided by Sauer and Geiger [22]. To that end, we will address claim HQC by performing a qualitative analysis on both MNIST and ImageNet. To verify claim IBR, we perform a loss ablation study in which we disable one loss at a time. Lastly, to address the main claim of the paper, namely ODR, we conduct a number of experiments on both MNIST and ImageNet to evaluate both out-of-distribution performance and spurious signal invariance of the invariant classifiers.

To provide further evidence to support claim ODR, we conduct additional experiments to visually explain the decisions made by the invariant classifiers based on gradient-based localization. For this purpose, we use a PyTorch implementation of GradCAM [11, 25], a class activation map method that weighs the 2D activations by the average gradient [25]. This method allows us to visualize the salient features on which the invariant classifiers base their predictions.

## 3.4 Computational requirements

We perform all experiments on a cluster whose nodes are equipped with Nvidia GeForce GTX 1080 Ti GPUs. Due to constraints in resources, we run most experiments once. As such, our experiments are indicative and not conclusive. Our reproducibility study comes at a total computational cost of 112 GPU hours (see Appendix D for more details).

# 4 EXPERIMENTAL RESULTS

## 4.1 Reproducibility study

**Evaluating counterfactual samples** To verify claim HQC, we qualitatively evaluate counterfactual (CF) samples generated using CGN models on each dataset. For all our reproducibility experiments, we use the available pretrained weights for CGN to generate CFs. We found inconsistencies while training the CGN from scratch and refer the reader to Section 4.2.1 for a deeper investigation. For both MNIST and ImageNet, our results indicate that the quality of the generated CFs matches with the quality of those reported in the original paper, as shown in Figure 1 and Figure 2 respectively. For ImageNet, although we can easily recognize the FoVs in the generated CFs, they are highly unrealistic.

<sup>1</sup>This variant of MNIST is generated by the authors themselves and can be generated using their repository.

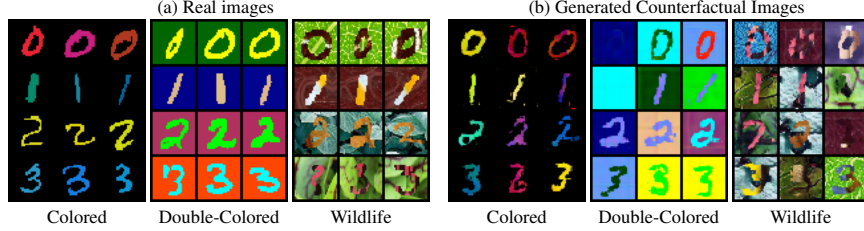


Figure 1: **Qualitative Analysis MNIST.** Left: Samples drawn from the different MNIST variations. Right: Counterfactuals generated by the CGN on MNIST variants. Notice that the CGN generates varying shapes, colors, and textures.

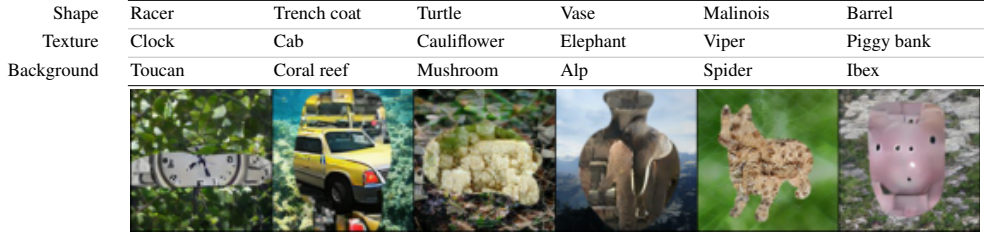


Figure 2: **Qualitative analysis ImageNet.** Counterfactuals generated by the CGN on ImageNet.

146 **Evaluating loss ablation** We attempt to reproduce the loss ablation study to verify claim IBR. The authors claim that  
 147 a CGN can be trained from scratch within 12 hours on a GTX 1080Ti GPU. However, when running the experiments  
 148 as described by the authors, the estimated training time exceeded 200 hours. Upon further inspection, we found an  
 149 alternative configuration file containing the hyperparameters the authors used to train the CGN that was inconsistent  
 150 with the default hyperparameters. Using these alternative hyperparameters, we managed to decrease the training time to  
 151 approximately 20 hours. While the inception score magnitude directly depends on the number of generated images  
 152 used for the calculation, the original paper did not specify the exact number of images used during the experiment. We  
 153 empirically found that using 2000 images provides inception scores that resemble those reported in the original paper.

154 The results in Table 2 indicate that the inception scores follow a similar trend as reported by the authors (marked as  
 155  $\times$ ). However, when disabling the texture loss, we found  $\mu_{mask}$  to be 0.4, whereas the original paper reported a value  
 156 of 0.9. This is a crucial difference, because the value of 0.9 of the original paper indicates a mask collapse, which  
 157 the authors use to support claim IBR. Nonetheless, we were able to support this claim by performing an additional  
 158 qualitative experiment. Specifically, if we look at some samples as shown in Appendix E, it is clear that the generated  
 159 texture still includes some background. This indicates that the independent mechanisms for texture and background are  
 160 no longer disentangled, which shows that the texture loss is indeed necessary.

161 **Evaluating invariant classifiers** We perform a number of experiments to verify claim ODR. Specifically, we  
 162 quantitatively evaluate the extent to which invariance is encoded in classifiers trained on CF data against those trained on  
 163 original data. We also evaluate classifiers trained on vanilla GAN-generated data as a baseline. Since the vanilla GAN  
 164 implementation was not provided in the released code, we implement it ourselves and refer the reader to Appendix F  
 165 for details and generated samples.

Table 2: **Loss Ablation Study.** We turn off one loss at the time.

$\mathcal{L}_{shape}$	$\mathcal{L}_{text}$	$\mathcal{L}_{bg}$	$\mathcal{L}_{rec}$	IS $\uparrow$	$\mu_{mask}$
$\times$	$\checkmark$	$\checkmark$	$\checkmark$	100.8   85.9	0.3   0.2
$\checkmark$	$\times$	$\checkmark$	$\checkmark$	186.5   198.4	0.4   0.9
$\checkmark$	$\checkmark$	$\times$	$\checkmark$	200.9   195.6	0.1   0.1
$\checkmark$	$\checkmark$	$\checkmark$	$\times$	19.3   38.4	0.4   0.3
$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	156.1   130.2	0.3   0.3
BigGAN (Upper Bound)				202.9	-

Table 3: **MNIST classification.** In the test-set, the texture and background are randomized; only the digits shape corresponds to the class.

Setting	C-MNIST		DC-MNIST		W-MNIST	
	Train $\uparrow$	Test $\uparrow$	Train $\uparrow$	Test $\uparrow$	Train $\uparrow$	Test $\uparrow$
O(riiginal)	99.7   99.5	37.6   35.9	100   100	10.5   10.3	100   100	10.8   10.1
GAN	99.6   99.8	32.5   40.7	100   100	10.6   10.8	99.9   100	11.2   10.4
CGN	99.4   99.7	92.3   95.1	94.8   97.4	86.5   89.0	95.5   99.2	81.4   85.7
O + GAN	99.6   99.8	41.5   40.7	100   100	10.0   10.8	100   100	11.1   10.4
O + CGN	99.2   99.7	95.9   95.1	96.9   97.4	85.5   89.0	96.8   99.2	62.8   85.7

166 On MNIST variants, we identify an inconsistency in the experimental setup stated in the paper and code. The paper  
 167 seems to suggest using a combination of original and CF dataset, but the code only uses CF data. As reported in Table 3,  
 168 we experiment with both and observe similar results for C-MNIST and DC-MNIST. Surprisingly, for CGN, adding  
 169 original data hurts the performance for W-MNIST (62.9 vs. 81.4). Apart from that, the majority of our results are within  
 170 5% variation from those reported in the paper (marked as  $\times$ ), which supports the broader claim of better generalization  
 171 even in the presence of spurious correlations (e.g., texture in case of colored MNIST).

172 To evaluate the invariance in classifier heads on IN-mini, we first reproduce the experiment regarding shape bias from  
 173 the original paper. The shape bias is defined as the fraction of test samples for which the predicted label matches the  
 174 shape label of the input image [8]. In this case, we evaluate labels with predictions from each head. As reported in  
 175 Table 4, our results are smaller in comparison to the IN-1k results reported in the original paper. Nonetheless, the  
 176 overall trend does support claim ODR. Additionally, we replicate the experiment regarding the evaluation of background  
 177 robustness. The paper uses the notion of BG-gap that measures classifiers’ reliance on background signal [28]. Our  
 178 results, shown in Table 5, again slightly deviate from the original paper but the trend supports claim ODR.

Table 4: **Shape vs. texture.** Evaluation of shape biases of independent classifiers.

Trained on	Shape Bias	top-1 $\uparrow$	top-5 $\uparrow$
IN + GCN/Shape	54.8		
IN + GCN/Text	16.7	74.0	91.7
IN + GCN/Bg	22.9		
IN-mini + GCN/Shape	49.1		
IN-mini + GCN/Text	20.5	56.2	79.1
IN-mini + GCN/Bg	25.7		

Table 5: **Backgrounds Challenge.** Evaluation of robustness against adversarially chosen backgrounds.

Trained on	IN-9 $\uparrow$	Mixed-Same $\uparrow$	Mixed-Rand $\uparrow$	BG-Gap $\downarrow$
IN	95.6	86.2	78.9	7.3
SIN	89.2	73.1	63.7	9.4
IN + SIN	94.7	85.9	78.5	7.4
Mixed-Rand	73.3	71.5	71.3	0.2
IN + CGN	94.2	83.4	80.1	3.3
IN-mini + CGN	86.8	73.2	68.3	4.9

179 To evaluate the effect of using more counterfactual datapoints or generating more counterfactual images per sampled  
 180 noise, Sauer and Geiger performed an MNIST Ablation Study in the original paper. Our reproduction for this experiment,  
 181 along with a more detailed description regarding the experiment and results, can be found in Appendix G.

## 182 4.2 Results beyond original paper

### 183 4.2.1 Improving CGN training on MNISTs

184 While training the CGN on the MNIST, we encountered an issue that was not mentioned in the original paper. During  
 185 the training process, we observed that the digit masks had a tendency of collapsing to an erroneous state, from where  
 186 the digits would no longer improve during training. For this reason, it was not possible for us to reproduce the CGN  
 187 training on the MNIST data using the default configuration. Therefore, we have proposed a solution that makes the  
 188 CGN training on the MNIST datasets more consistent. Details regarding our solution can be found in Appendix C.

### 189 4.2.2 Explainability analysis for invariant classifiers

190 While the reproduced experiments for the original paper provide some support for claim ODR, these results primarily  
 191 show the effect of using counterfactuals on test accuracy performance. However, it is not directly clear from these  
 192 quantitative experiments if the performance increase is actually due to the fact that the use of counterfactuals ensures  
 193 that the classifier focuses on the right correlations (e.g., shape) and not spurious ones (e.g., background). To further  
 194 verify the validity of claim ODR, we provide two additional analyses that combine qualitative and quantitative measures  
 195 to evaluate the behaviour of the counterfactual classifiers.

196 **What does the latent feature space look like?** First, we visualize learnt classifier features using t-SNE for a subset  
 197 of the test set of original and counterfactual (CF) data for C-MNIST. Figure 3(a) shows that a classifier trained on CF  
 198 data is indeed invariant to spurious correlations (e.g. digit color). Figure 3(b) shows that a classifier trained on CF data  
 199 is also better at representing OOD samples (e.g. counterfactuals). Interestingly, the latter figure also shows that the  
 200 CF-trained classifier tends to group the clusters for 4-7-9 and 3-5-8 close to each other, which was not the case for the  
 201 classifier trained on original data. These digits are also close in shape in reality, which suggests that the model is rightly  
 202 focusing on the shape while ignoring texture. The results for other MNIST variants are consistent with this finding.

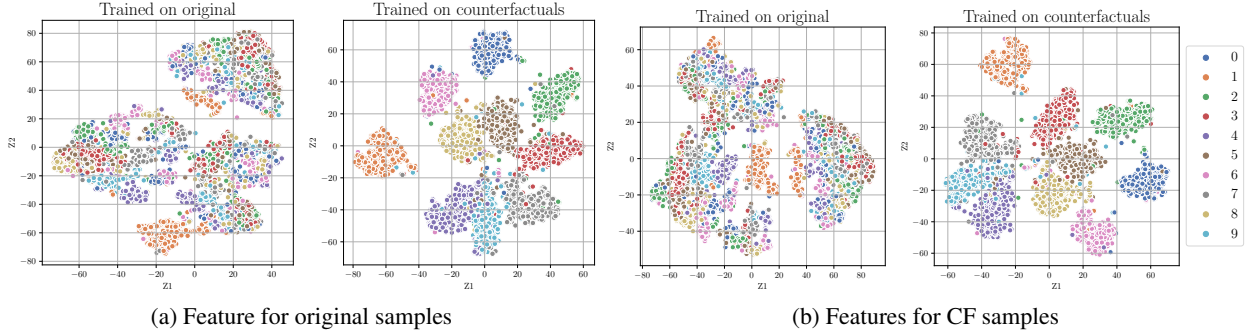


Figure 3: **Feature visualization.** Feature space of a CNN classifier trained on original/CF data for colored MNIST.

203 **What features does the model focus on?** Second, we perform an experiment to visualize a spatial heatmap of areas  
 204 that the model focuses on to make a prediction. Based on claims ODR and IBR, we would expect the different heads to  
 205 operate separately from one another, while being completely invariant to the other FoVs. In order to generate the spatial  
 206 heatmaps we use GradCAM. Some qualitative samples are shown in Figure 4. In addition to the qualitative analyses,  
 207 using GradCAM provides the opportunity to formulate another quantitative measure to validate claims ODR and IBR.  
 208 This quantitative analysis aims to measure if CF-trained models focus on shape more than those trained on original data.  
 209 To this end, we compute the mean Intersection of Union (IoU) between GradCAM heatmaps and binarized digit masks  
 210 on the test set. We note that a classifier trained on CF data is consistently outperforms the classifier on original data.

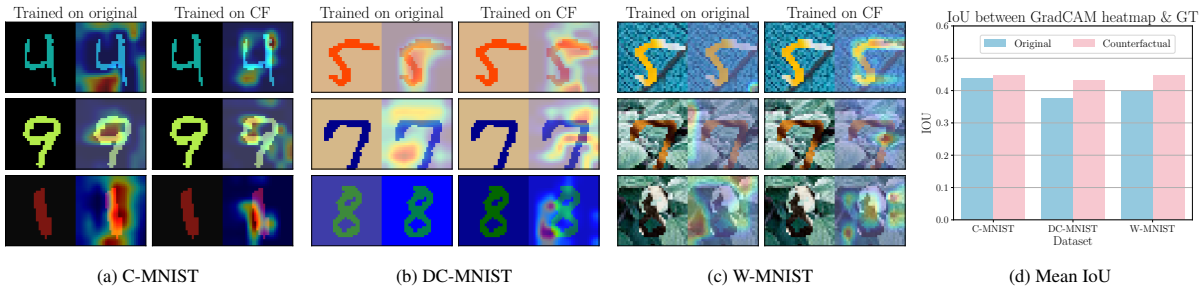


Figure 4: **Explainability analysis:** (a) to (c): Visualization of GradCAM heatmaps for samples from each of the MNIST datasets. (d): Mean IoU between GradCAM heatmaps and ground truth binarized digit masks.

211 While the quantitative results using the IoU metric cannot be performed on the ImageNet data, due to the lack of ground  
 212 truth binary object maps, it is possible to evaluate the qualitative performance of the independent mechanisms using  
 213 GradCAM. As shown in Appendix H, the individual classifier heads tend to focus on meaningful aspects.

#### 214 4.2.3 OOD generalization for invariant classifiers

215 In order to provide further evidence for the claim ODR, we test the model performance on alternative ImageNet datasets,  
 216 which are designed to evaluate out-of-distribution robustness. Specifically, we evaluate the performance on ImageNet-A  
 217 (natural adversarial examples) [13], ImageNet-Sketch [27] and Stylized-ImageNet [9], and compare with a ResNet-50  
 218 baseline that is pretrained on IN-1k. Surprisingly, we find that the finetuned CGN-based ensemble performs worse on  
 219 all specified OOD-benchmarks, compared to the pretrained ResNet-50 baseline as shown in Table 6.

Table 6: Comparison of top-1 accuracy of invariant classifier with pretrained ResNet on OOD benchmarks.

Model	Pretrained	Finetuned	IN-mini $\uparrow$	IN-A $\uparrow$	IN-Sketch $\uparrow$	IN-Stylized $\uparrow$
ResNet-50	IN-1k	-	75.580	3.400	24.092	19.218
CGN Ensemble	IN-1k	IN-mini + CF	56.793	1.387	11.775	17.188

## 220 5 DISCUSSION

221 Throughout this work, we have conducted several experiments to reproduce the main results from the research by  
222 Sauer and Geiger [22]. The results of our reproducibility study provide support for their claims, as we were largely  
223 able to reproduce the original results. Specifically, our results showed that the test accuracy for the MNIST classifiers  
224 greatly improved when using generated counterfactual datasets. Then, we were able to use the ImageNet-mini dataset  
225 to achieve similar performance trends compared to the original paper in terms of shape versus texture bias evaluation,  
226 and the background robustness evaluation. However, based on the qualitative analyses for claim HQC, it is clear that  
227 the quality of the generated counterfactual images could still be improved. Specifically, we have observed some distinct  
228 failure cases regarding the quality of generated counterfactual images, which are described in Appendix I.

229 Interestingly, while the loss ablation study provided similar results to what the authors reported in the original paper,  
230 we did obtain different results for the experimental run without texture loss. As the authors used this study to provide  
231 evidence for claim IBR, this difference is quite significant. Nonetheless, qualitative analysis of the images that were  
232 generated without texture loss revealed that the quality of the generated images indeed reduced when the texture loss  
233 was omitted. Although this does provide support for claim IBR, it also shows that the IS and  $\mu_{mask}$  metrics used by  
234 the authors in the loss ablation study may not be sufficient to support their claims. Since the loss ablation study is  
235 therefore not conclusive, further research is required to investigate if the inductive biases introduced by the authors are  
236 indeed ‘appropriate’. The results from our additional experiments provide further evidence that counterfactual images  
237 generated with the proposed CGN architecture can be used to train classifiers that are more robust against spurious  
238 signals. Using GradCAM, we were able to visualize this behaviour and formulate a quantitative performance metric.

239 Overall, the experiments from the original paper were largely reproducible, and their main claims seem reasonably sub-  
240 stantiated but could benefit from additional evidence in future research. The code implementation of our reproducibility  
241 study is publicly available <sup>1</sup>.

242 **Limitations** Unfortunately, we did encounter some difficulties during the reproduction process. First, since our model  
243 was trained on IN-mini, we were not able to reproduce the exact same results as the original paper. However, despite  
244 the slightly deviating results, the overall trends in the results seem to correspond well with the original results. Second,  
245 as some experimental setup information was missing from the original paper, we had to rely on the default parameter  
246 configuration files that were provided in the original code implementation, even though we can not be completely  
247 certain that these parameters were used for the original experiments.

### 248 5.1 Reflection: What was easy, and what was difficult?

249 The original paper provides an extensive appendix with implementation details and hyperparameters. Beyond that, the  
250 original code implementation was publicly accessible and well structured. As such, getting started with the experiments  
251 proved to be quite straightforward. The implementation included configuration files, download scripts for the pretrained  
252 weights and datasets, and clear instructions on how to get started with the framework.

253 Nonetheless, reproducing the original results turned out to be far from trivial as the setup of some of the experiments  
254 required severe modifications to the provided code. Additionally, some details required for the implementation are not  
255 specified in the paper or inconsistent with the specifications in the code (e.g., the GAN as mentioned in Section 3).  
256 Lastly, in evaluating robustness to OOD, getting the baseline model to work and obtaining numbers similar to those  
257 reported in the respective papers was challenging, partly due to baseline model inconsistencies within the literature.

### 258 5.2 Communication with original authors

259 We have reached out to the original authors to get clarifications regarding the setup of some of the experiments. For  
260 example, we asked the authors if they could share pretrained weights from the classifiers that were trained on full  
261 ImageNet, and which type of GAN architecture was used for the MNIST experiments. Unfortunately, we received a  
262 late response and only a subset of our questions was answered, and as a result we were not able to fully verify whether  
263 our design choices were consistent with those of the original paper.

---

<sup>1</sup><https://github.com/anonymous-user-256/mlrc-cgn>



## References

- [1] M. A. Alcorn, Q. Li, Z. Gong, C. Wang, L. Mai, W.-S. Ku, and A. Nguyen. Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects, 2019.
- [2] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. Invariant risk minimization, 2020.
- [3] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020.
- [4] A. Brock, J. Donahue, and K. Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [5] A. J. DeGrave, J. D. Janizek, and S.-I. Lee. Ai for radiographic covid-19 detection selects shortcuts over signal. *Nature Machine Intelligence*, pages 1–10, 2021.
- [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [7] I. Figotin. Imagenet 1000 (mini). <https://www.kaggle.com/ifigotin/imagenetmini-1000>, 2019.
- [8] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *CoRR*, abs/1811.12231, 2018. URL <http://arxiv.org/abs/1811.12231>.
- [9] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bygh9j09KX>.
- [10] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- [11] J. Gildenblat and contributors. Pytorch library for cam methods. <https://github.com/jacobgil/pytorch-grad-cam>, 2021.
- [12] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, D. Pedreschi, and F. Giannotti. A survey of methods for explaining black box models, 2018.
- [13] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song. Natural adversarial examples. *CVPR*, 2021.
- [14] N. Inkawhich. Dcgan tutorial. [https://pytorch.org/tutorials/beginner/dcgan\\_faces\\_tutorial.html](https://pytorch.org/tutorials/beginner/dcgan_faces_tutorial.html), 2021.
- [15] D. Kaushik, E. Hovy, and Z. C. Lipton. Learning the difference that makes a difference with counterfactually-augmented data. *arXiv preprint arXiv:1909.12434*, 2019.
- [16] O. Lang, Y. Gandelsman, M. Yarom, Y. Wald, G. Elidan, A. Hassidim, W. T. Freeman, P. Isola, A. Globerson, M. Irani, et al. Explaining in style: Training a gan to explain a classifier in stylespace. *arXiv preprint arXiv:2104.13369*, 2021.
- [17] Y. Ming, H. Yin, and Y. Li. On the impact of spurious correlation for out-of-distribution detection, 2021.
- [18] M. L. Olson, R. Khanna, L. Neal, F. Li, and W.-K. Wong. Counterfactual state explanations for reinforcement learning agents via generative deep learning. *Artificial Intelligence*, 295:103455, 2021.
- [19] J. Pearl. *Causality*. Cambridge university press, 2009.
- [20] X. Qin, Z. Zhang, C. Huang, M. Dehghan, O. R. Zaiane, and M. Jagersand. U2-net: Going deeper with nested u-structure for salient object detection. *Pattern Recognition*, 106:107404, 2020.
- [21] A. Rosenfeld, R. Zemel, and J. K. Tsotsos. The elephant in the room, 2018.

- 305 [22] A. Sauer and A. Geiger. Counterfactual generative networks. In *International Conference on Learning Representations (ICLR)*, 2021.  
306
- 307 [23] A. Sauer and A. Geiger. Counterfactual generative networks github. [https://github.com/](https://github.com/autonomousvision/counterfactual_generative_networks)  
308 [autonomousvision/counterfactual\\_generative\\_networks](https://github.com/autonomousvision/counterfactual_generative_networks), 2021.
- 309 [24] B. Schölkopf. Causality for machine learning. *arXiv preprint arXiv:1911.10500*, 2019.
- 310 [25] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2):  
311 336–359, Oct 2019. ISSN 1573-1405. doi: 10.1007/s11263-019-01228-7. URL [http://dx.doi.org/10.](http://dx.doi.org/10.1007/s11263-019-01228-7)  
312 [1007/s11263-019-01228-7](http://dx.doi.org/10.1007/s11263-019-01228-7).  
313
- 314 [26] S. Verma, J. Dickerson, and K. Hines. Counterfactual explanations for machine learning: A review. *arXiv preprint arXiv:2010.10596*, 2020.  
315
- 316 [27] H. Wang, S. Ge, Z. Lipton, and E. P. Xing. Learning robust global representations by penalizing local predictive  
317 power. In *Advances in Neural Information Processing Systems*, pages 10506–10518, 2019.
- 318 [28] K. Xiao, L. Engstrom, A. Ilyas, and A. Madry. Noise or signal: The role of image backgrounds in object  
319 recognition, 2020.
- 320 [29] K. Xiao, L. Engstrom, A. Ilyas, and A. Madry. Noise or signal: The role of image backgrounds in object  
321 recognition. *ArXiv preprint arXiv:2006.09994*, 2020.

## 322 A Counterfactual Generative Network Architecture

323 In Figure 5, we provide an overview of the architecture of the CGN as provided in the paper. It illustrates how the  
 324 CGN is split into four mechanism: the shape mechanism, the texture mechanism, the background mechanism, and  
 325 the composer. Each mechanism takes a noise vector  $\mathbf{u}$  and a label  $y$  as input. To generate a counterfactual image, we  
 326 sample  $\mathbf{u}$  and then sample a separate  $y$  for each mechanism Sauer and Geiger [22].

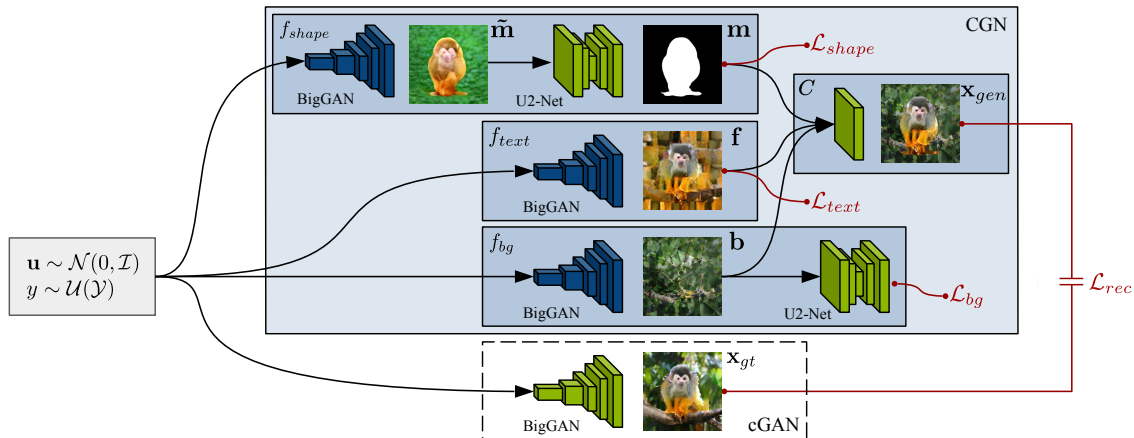


Figure 5: **CGN architecture.** Components with trainable parameters are blue, components with fixed parameters are green [22]. The dotted lines indicate that the cGAN is only used for training [22].

## 327 B Counterfactual images and explainability in artificial intelligence

328 One of the primary contributions of the work by Sauer and Geiger [22] is the proposed method to create high-quality  
 329 ‘counterfactual’ images, which can be used to make a classifier more robust to spurious signals. As the concept of  
 330 *counterfactual explanations* is closely related to the idea of explainable artificial intelligence (XAI) but is never explicitly  
 331 mentioned in the paper, we first want to place the article in a broader context to achieve a deeper understanding of how  
 332 the considered work relates to other developments within this field of research [3].

333 Based on the review by Verma et al. [26], approaches for explainability in machine learning can be roughly divided  
 334 into one of two categories: (i) methods that use inherently interpretable and transparent models, and (ii) methods that  
 335 generate post-hoc explanations for opaque models. The idea of counterfactual explanations belongs to the example-  
 336 based approaches within the category of post-hoc explanations, that seek to offer explanations by either providing  
 337 datapoints that receive the same prediction label as the observed datapoint, or by providing datapoints whose prediction  
 338 label is different from the observed datapoint.

339 Consider the example where a classifier is trained to distinguish images from polar bears and American black bears.  
 340 Given an image that has been classified by the model as a black bear, we could attempt to provide a post-hoc explanation  
 341 for the model’s prediction using a visual counterfactual explanation (i.e., a modified version of the input image that  
 342 would be classified as a polar bear instead). These explanations can, for example, be generated using techniques such as  
 343 StyleEx [16]. A reasonable visual counterfactual explanation could consist of the input image, modified such that the fur  
 344 of the black bear is now colored white. However, as most images of polar bears have a snow-background, and most  
 345 images of American black bears likely do not, it is possible that the suggested visual counterfactual explanation still  
 346 contains a black bear, but now on a snowy background.

347 In this case, one could argue that the background-explanation that is captured by the model is a spurious signal. That is,  
 348 the classifier ‘falsely’ makes predictions on the background, even though the background, in reality, does not affect  
 349 the actual object itself. Although this spurious signal might seem innocent within the context of this example, other  
 350 spurious signals can play a role in a variety of high stake deep learning applications, such as AI in medical-imaging  
 351 [5] and networks trained for military purposes [12]. While counterfactual explanations are thus capable of *revealing*  
 352 such spurious signals, the proposed method using counterfactual images by Sauer and Geiger provides an approach to  
 353 *mitigate* this effect.

## 354 C Improved CGN Training for MNIST

355 While training the CGN on the MNIST, we encountered an issue that was not mentioned in the original paper. During  
 356 the training process, we observed that while some digits were captured almost perfectly by the model, other digit masks  
 357 seemed to collapse to a state where there was a black circular shape in the center of the image with a surrounding  
 358 white border (see Figure 6). When using the generated counterfactual datasets from these imperfect models to train  
 359 a classifier, we then observed that the number of ‘correct’ (i.e., non-collapsed) images correlated strongly with the  
 360 classifier performance.

361 Any attempt to remedy this issue using adjusted hyperparameter configurations proved to be ineffective, because the  
 362 hyperparameter names in the provided default configuration-files did not directly correspond to the descriptions given  
 363 in the original paper. This observation inspired a solution where we add an extra loss term to the training objective,  
 364 which penalizes mask-pixels at the borders of the image. Specifically, if we define the edge region  $\mathcal{E}$  as the set of pixels  
 365 that are within  $s$  pixels from the edge, the edge loss function can be defined as the sum of all pixel values  $m_i$  within the  
 366 specified edge region:

$$\mathcal{L}_{edge}(\mathbf{m}) = \mathbb{E}_{p(\mathbf{u}, \mathbf{y})} \left[ \frac{1}{N} \sum_{i=1}^N m_i \cdot [i \in \mathcal{E}] \right], \quad (3)$$

367 where  $N$  denotes the number of pixels in mask  $\mathbf{m}$ , and  $[\cdot]$  denotes the Iverson bracket. As the original MNIST images  
 368 in the training and test datasets often contain almost no pixels at the borders, this loss function returns values close to 0  
 369 for all ground truth MNIST images. During our experiments, we used a border size of 3 pixels, as this configuration  
 370 seems to perform well to mitigate the mask-collapse issue, while still giving loss values close to 0 for the original  
 371 MNIST images. By using this extra loss function, the training process became much more consistent and lead to an  
 372 average classifier test accuracy of 89.8% for the Colored MNIST dataset, which is close to what was reported in the  
 373 original paper.

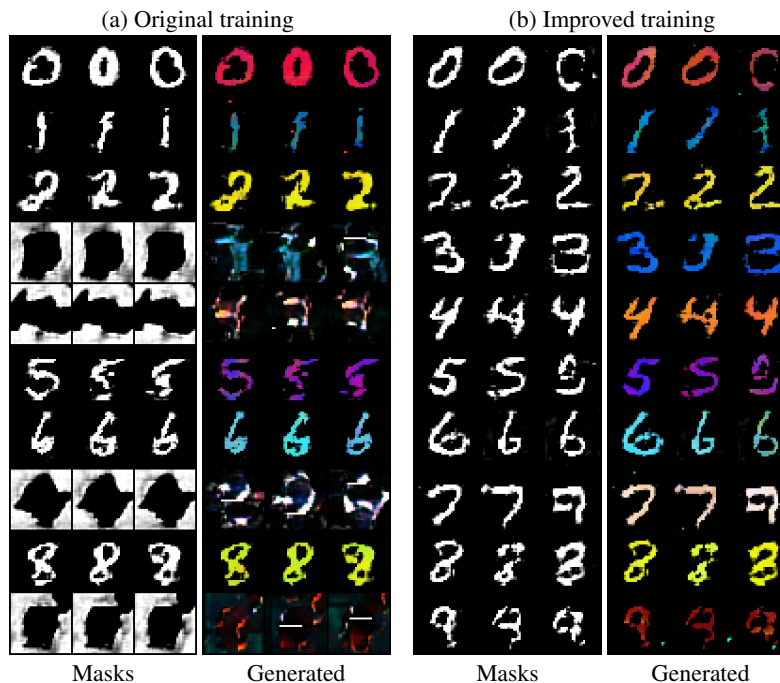


Figure 6: **Qualitative edge loss evaluation.** Adding the edge loss significantly improves CGN training on colored MNIST.

374 In Figure 10, we show that our modified training formulation improves the quality of generated images. In particular,  
 375 we notice that incorporating  $\mathcal{L}_{edge}$  in the mask loss, on average, noticeably decreases the number of non-broken images.

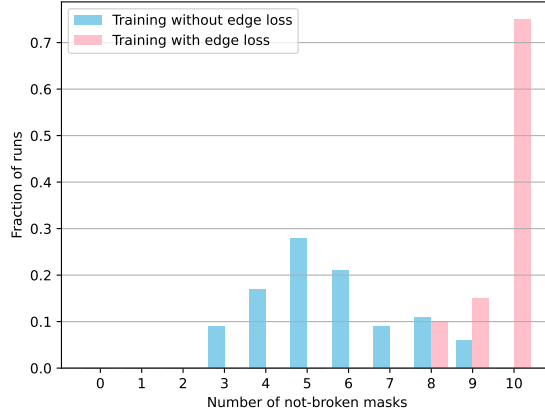


Figure 7: **Quantitative edge loss evaluation.** The fraction of experiment runs for each number of ‘correct’ digits.

### 376 D Computational Cost Taxonomy

Table 7: **Cost taxonomy.** Overview of the computational cost associated with each experiment.

Experiment type	Experiment name	Support of Claim	Section	Computational Cost (GPU Hours)
Reproducibility Study	Evaluating counterfactual samples	HQC	4.1	0.0
	Required Inductive Biases	IBR	4.1	84.0
	Evaluating invariant classifiers: MNIST	ODR	4.1	6.0
	Evaluating invariant classifiers: IN-Mini	ODR	4.1	8.0
	Ablation study (Appendix G)	ODR	4.1	14.0
Additional results	Improved CGN Training	HQC	4.2.1	48.0
	Explainability analysis: MNIST	ODR	4.2.2	< 1.0
	Explainability analysis: IN-Mini	ODR	4.2.2	< 1.0
	OOD generalization evaluation	ODR	4.2.3	< 1.0

### 377 E Qualitative Analysis of Loss Ablation Study

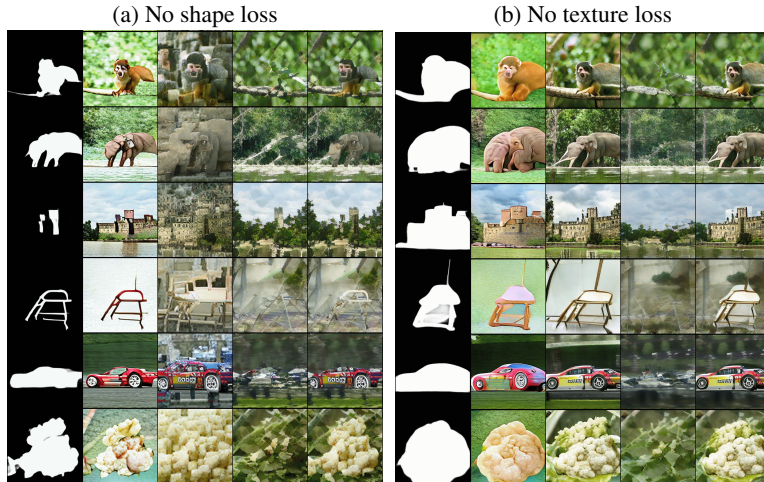


Figure 8: **Qualitative Loss Ablation.** Comparison between IM outputs when excluding the shape loss and texture loss. From left to right:  $m$ ,  $\tilde{m}$ ,  $f$ ,  $b$ ,  $x_{gen}$  as described in Section 2.

378 **F GAN-based Baseline for MNISTs**

379 We follow the ConvNet-based architecture for the generator inspired by PyTorch DCGAN tutorial and retain the linear  
 380 discriminator as is used by Sauer and Geiger [22]. We only use binary cross entropy loss for adversarial training of both  
 381 G and D. All necessary hyperparameters are same as for the CGN training. These along with pretrained weights can be  
 382 found in our code repository.

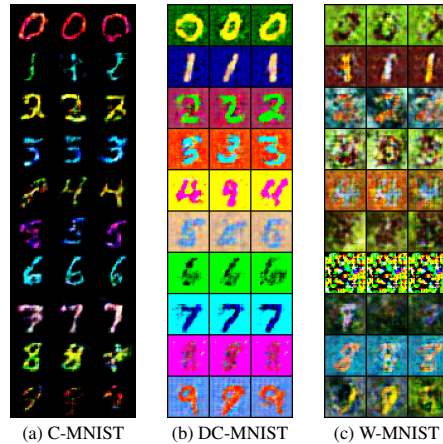


Figure 9: **GAN samples.** Samples generated by a GAN baseline on MNIST variants.

383 **G Reproduced MNIST Ablation Study**

384 Figure 10 shows our reproduced results for the MNIST ablation study. Our results show that using more counterfactual  
 385 datapoints generally improves the test accuracy, although this was not the case for the Colored MNIST dataset, where  
 386 the test accuracy decreased when using  $10^6$  counterfactual datapoints instead of  $10^5$ . However, the difference in  
 387 performance is only minor. The differences in CF ratios do not seem to have a significant effect on the test accuracies.  
 388 These results seem to support the claim from the original paper that using more counterfactual images always increases  
 389 the test domain results for MNIST datasets, although there only seems to be a significant performance increase when  
 390 using  $10^5$  datapoints instead of  $10^4$ . Using even more datapoints does not seem to provide a significant increase in  
 391 performance.

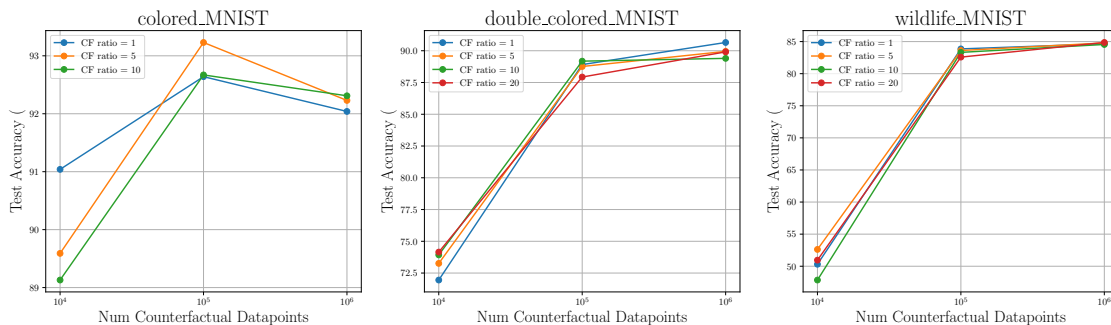


Figure 10: **MNIST ablation study.** We evaluate the impact of using more counterfactual data and generating more counterfactuals per sampled noise on the measured test accuracy.

392 **H GradCAM samples on ImageNet-mini**

393 A classifier trained jointly on original and CF data is expected to have encoded invariances for certain attributes and  
 394 distinctiveness for others. Recall that the proposed classifier architecture for ImageNet is an ensemble with three heads  
 395 for shape, texture and background. We pose the question: What spatial aspects of an image does each head *focus* on and



396 what prediction does it lead to? We answer this qualitatively by analyzing GradCAM heatmaps for outputs of each of  
 397 the heads as well as the averaged ensemble output. In general, the individual heads tend to focus on meaningful aspects,  
 398 as shown in Figure 11, background head focuses on background. Further, for original images, we observe that a correct  
 399 prediction often relies on shape (e.g., *puck* in Figure 11a) or texture (e.g., *goldfinch*). In some cases, it correctly relies  
 400 on background (e.g., *castle*). For counterfactuals, surprisingly, in most cases we found that the label predicted from  
 401 shape, although correct, is dominated by incorrect label from background and texture. This may be a symptom of either  
 402 insufficient counterfactual training data or the use of IN-mini instead of IN-1k. We further note that texture often drives  
 403 the label decision for counterfactuals.

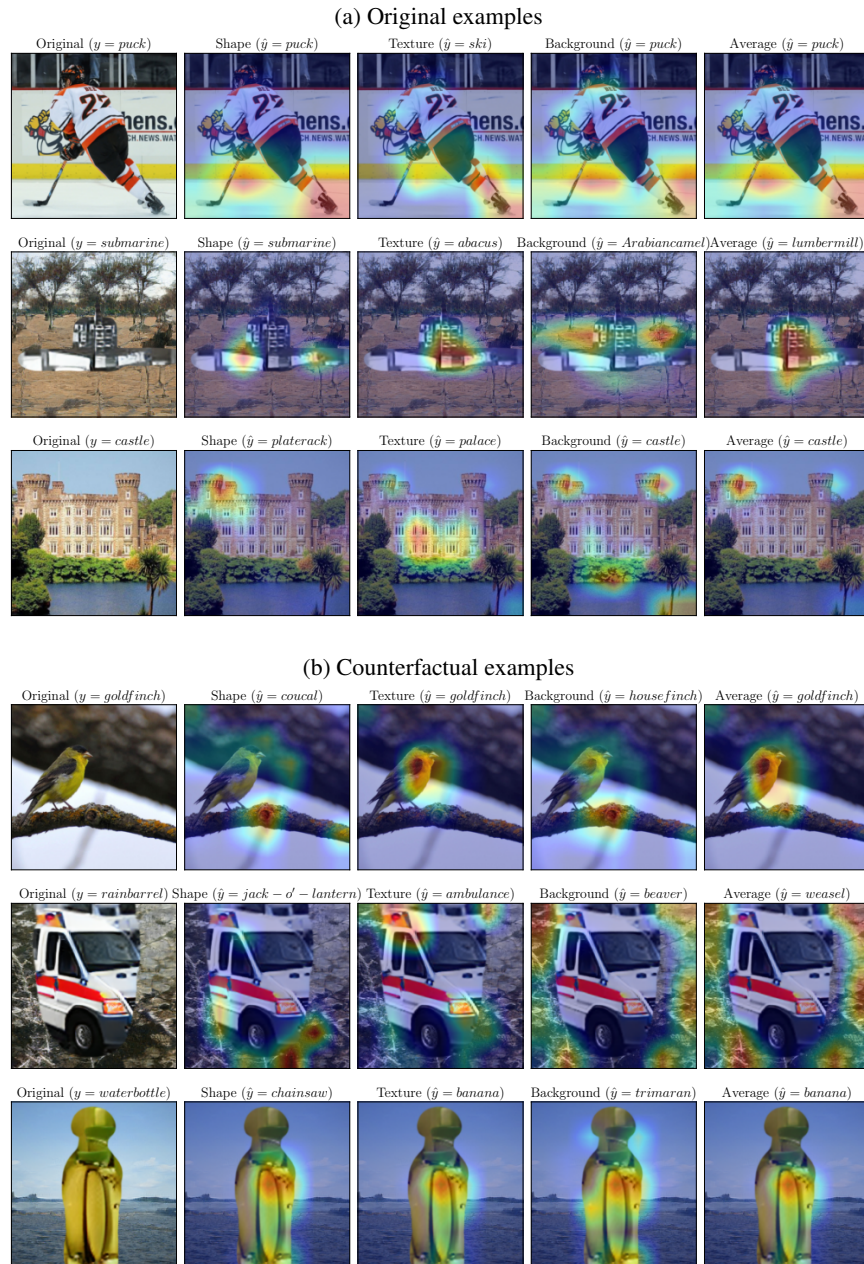


Figure 11: **Explainability Analysis ImageNet.** GradCAM heatmaps visualized with respect to individual head outputs for original and counterfactual samples. The corresponding ground truth labels and predictions are provided too.

404 **I Some failure modes in CGN-generated samples**

405 Since generation of high-quality counterfactuals is one of the main claims of the paper, we perform a deeper qualitative  
406 analysis to observe if there exist typical failure modes. Based on anecdotal evidence, we note the following observations.

407 **Texture-background entanglement for small objects** For cases with small objects on a uniform background, such  
408 as the bird kite in sky, shown in Figure 12(a), or skiing on snow, shown in Figure 12(b), we see consistent  
409 entanglement between texture and background.

410 **Objects with complex texture** We observe that objects with complicated texture, such as crossword puzzle,  
411 shown in Figure 12(c), result in poorly recovered texture by the CGN.

412 **Complex scenes** As one would expect, the CGN approach does not generalize to complex scenes since it assumes a  
413 simplistic causal structure. We show an example of this in Figure 12(d).



Figure 12: **Failure modes.** Cases highlighting some common failure modes in samples generated using CGN.