
Reproducibility study - Counterfactual Generative Networks

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1
2 **Scope of Reproducibility** In this study, we worked on the reproducibility of the results in the paper Counterfactual
3 Generative Networks by Axel Sauer, Andreas Geiger. The study is performed based on the following claims;

- 4 • Counterfactual generative network (CGN) can generate high-quality counterfactual images with direct control
5 over shape, texture, and background.
- 6 • Using generated counterfactual images in training data set improves the classifier's out-of-domain robustness.
- 7 • Using generated counterfactual images in the training data set only marginally degrades overall accuracy.

8 **Methodology** Source code used in the original paper was already provided by the authors and implemented in Pytorch.
9 Code was adapted for different experimentation purposes. Additionally, authors used some pre-trained networks in their
10 experiments. Original paper includes a link to these networks' implementation as well.

11 **Results** We managed to reproduce most of the results in the original paper. We had some difficulties reproducing the
12 first claim, but the results of our experiments support the second and the third claim.

13 **What was easy** The architecture of the networks was explained clearly in the paper and it was relatively easy to
14 comprehend. Implementation-wise, the code was clean enough to run without requiring an extensive debugging.
15 Appendix in the paper provided quite many visualization and detailed explanation regarding the experiments, including
16 the failed cases. This gave us an insight about the limitations in the models' performance.

17 **What was difficult** Main difficulty in the experiments was that the computation time required for the model training with
18 ImageNet data set. It is approximated to take about 214 hours to conduct a single experiment, while running on a cluster
19 computing system. To complete the experiments in the given time frame, subset of the ImageNet (ImageNet-mini) is
20 used.

21 **Communication with original authors** We contacted with the authors regarding the discrepancies between the code
22 and the paper, and the unaligned results. Authors clarified our concerns several times in different occasions.

23 **1 Introduction**

24 Deep neural networks (DNNs) are the fundamental learning algorithms which are widely used in the field of machine
25 learning. Although the DNNs perform well in many tasks, they still struggle with handling the unseen data. Data bias is
26 seen as the main factor in the failed cases. If the model has always seen a particular object in a certain background,
27 it tends to correlate the object with the background, and the model fails to recognize the object when it appears in a
28 different background, which is a significant obstacle towards having a generalized model.

29 Data augmentation is seen as a strong regularizer and an efficient way to extend the training data set in machine
30 learning algorithms, by Wong et al. (2016). Previous studies by Goyal et al. (2018) showed that data augmentation with
31 synthetic images is a promising solution. Authors designed a new generative model called counterfactual generative
32 network (CGN), which generates counterfactual images based on two main assumptions; independent mechanism and
33 the composition mechanism. Independent mechanism suggests that the modules that generate the synthetic image are
34 independent of each other (dos Santos Tanaka and Aranha (2019)). This way spurious correlation can be minimized.

35 CGN generates an image by combining three independent components which are defined as shape, texture and
36 background, and those components are combined analytically, following a certain equation. In this paper, same
37 assumptions are accepted for the sake of assessing the reproducibility of the experiment results.

38 **2 Scope of reproducibility**

39 Focus of our work is to reproduce the general trends in the experimental results, such as an increase in the performance
40 of the classifier which is trained on counterfactual images in addition to the original data set. In this study, we worked
41 to validate the following claims in the original paper;

- 42 1. CGN can generate high-quality counterfactual images with direct control over shape, texture, and background.
- 43 2. Using generated counterfactuals in the training data set improves the classifier's out-of-domain robustness
- 44 3. Using generated counterfactuals in training data set only marginally degrades overall accuracy

45 Claim 1, which is supported by experiments found in Section 3.4.1 and Fig 4, 2, 9, 8 was not proven correct. Second
46 claim, which is assisted by experiments described in Section 3.4.2 and Table 2 was found correct. Finally, third claim is
47 supported by Section 3.4.3 and Table 1 and 5 and 4 and proven correct. In the following sections, methodology we
48 adopted in this study is explained. Following that, claims and the results are presented. Finally, we discuss the strengths
49 and the weaknesses of the original paper in the discussion section.

50 **3 Methodology**

51 **3.1 Implementation**

52 The base code ¹ is provided by the authors in the original paper. It was well documented and sufficiently clear to run
53 without requiring any debugging.

54 **3.2 Model descriptions**

55 To investigate claims made by the authors in the original papers, we carried out experiments using counterfactual
56 generative networks (CGNs). In this section we describe the architectures of the two different CGNs we used.

57 **MNIST CGN** : In the original paper, it is assumed that the generative process of the counterfactual images can be
58 decomposed into three independent mechanisms; shape, texture and background mechanism. For the MNIST data set,
59 mechanisms for texture and background are designed with the exact same structure, while shape mechanism has a
60 slightly different structure.

¹https://github.com/autonomousvision/counterfactual_generative_networks

²This image is taken from the original paper Sauer and Geiger (2021).

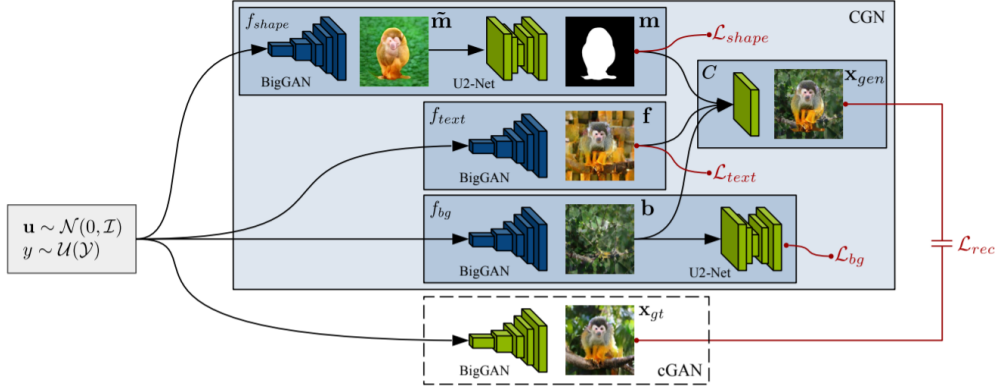


Figure 1: **Counterfactual Generative Network (CGN)** Here, the architecture of the CGN is displayed. The network consists of four main mechanisms. These are f_{shape} for shape component, f_{text} for texture, f_{bg} for background and C for composition. Pretrained models are shown in green, while the models with trainable parameters are shown in blue. Performance of the CGN is assessed via the reconstruction loss, which is computed using the conditional GAN (cGAN) outputs. In the model, cGAN took part only in the training process. Each mechanism receives a Gaussian noise vector \mathbf{u} and a label \mathcal{Y} . The loss values, which is shown in red, are minimized during the training process. Counterfactual images are generated using a noise vector and independently sampled labels (one label for each mechanism)².

61 In the generation of a new image, all three components are merged based on a certain composition mechanism, which is
 62 the second assumption in the original paper. Most important feature of the composition mechanism is that it is defined
 63 analytically, not learned by any model. Composition mechanism is defined as follows.

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

64 where the \mathbf{m} is the mask (shape component), \mathbf{f} is the foreground and \mathbf{b} is the background. \odot is used to represent the
 65 element-wise multiplication.

66 **ImageNet CGN** Architecture of the CGN designed for ImageNet is displayed in Fig 1. Independent mechanism and
 67 composition mechanism assumptions are applied in the ImageNet CGN as well. Additionally, several loss values are
 68 computed to take part in the model training. Since the mechanisms that comprise the generative model are independent
 69 of each other, each mechanism has its own loss value. Loss values can be found in red text in Figure 1. \mathcal{L}_{shape} ,
 70 \mathcal{L}_{text} , \mathcal{L}_{bg} correspond to the loss in shape, texture, background mechanisms, respectively, and \mathcal{L}_{rec} corresponds to
 71 the reconstruction loss which is computed using the output of the generative model and the output of the pre-trained
 72 BigGAN model. In the training process, all loss values are linearly combined and jointly optimized. Overall loss is
 73 calculated as follows.

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{rec} + \mathcal{L}_{shape} + \lambda_5 \mathcal{L}_{text} + \lambda_6 \mathcal{L}_{bg} \\ &= \lambda_1 \mathcal{L}_{L1} + \lambda_2 \mathcal{L}_{perc} + \lambda_3 \mathcal{L}_{binary} + \lambda_4 \mathcal{L}_{mask} + \lambda_5 \mathcal{L}_{text} + \lambda_6 \mathcal{L}_{bg} \end{aligned} \quad (2)$$

74 where $\lambda_1 = 100$, $\lambda_2 = 5$, $\lambda_3 = 300$, $\lambda_4 = 500$, $\lambda_5 = 5$, $\lambda_6 = 2000$. Authors did not provide any information regarding
 75 the calculation of the λ values.

76 Different than the MNIST dataset, trainable BigGAN models are placed in each independent mechanism. Furthermore,
 77 shape mechanism and background mechanism comprise a pre-trained U2-Net to process the output of the BigGAN
 78 models. Additionally, another pre-trained BigGAN model is used to generate non-counterfactual images using the given
 79 noise and the label, and it is mainly used to compute the reconstruction error and train the CGN.

80 3.3 Datasets

81 Two main datasets are used in the experiments; MNIST and ImageNet-1k. MNIST data set consists of three subsets;
 82 colored-MNIST, double-colored MNIST and the wildlife MNIST. Description of the datasets can be found below.

³This equation is provided by the original paper Sauer and Geiger (2021).

83 **Colored MNIST:** It has images 50k for training and 6k images for testing set. Ten different colors are selected and each
84 color is assigned to a single type of digit as their mean color in the training set. In the test set, the colors are randomly
85 assigned to the digits.

86 **Double-colored MNIST:** It has 50k for samples for training and 6k samples for testing. The main difference with the
87 colored MNIST is that the background is also colored with the one of the selected colors.

88 **Wildlife MNIST:** It has 50k for training and 6k for testing. Ten distinct texture images are selected from the striped
89 class which is provided by Cimpoi et al. (2013) for the texture. Ten other distinct texture images are selected from the
90 veiny class from the same source, for the background.

91 **ImageNet-mini:** It has 1k classes, 34.752 samples for training and 10k samples for validation.

92 3.4 Experimental setup

93 In this section we describe the experiments performed to investigate the claims described in section 2. For comprehensi-
94 bility, we list the claims and corresponding experiments. In this study, MNIST models are trained on a single RTX
95 1080 GPU located in a cluster computer system. A device with more memory, such as an Titan RTX, is advised when
96 training ImageNet models.

97 3.4.1 Claim 1

98 The following are experiments carried out to investigate the claim ‘A CGN can generate high-quality counterfactual
99 images with direct control over shape, texture, and background.’

100 **MNIST counterfactuals** We trained 3 CGN’s, each on one of the MNIST datasets (described in Section ??) using
101 code provided by the authors, which can be found and run [HERE](#). The architecture of the CGN’s trained on the MNIST
102 datasets is described in Section 3.2. During training we sample counterfactuals generated by the model and we compare
103 counterfactuals generated by the trained CGN to examples in Sauer and Geiger (2021). Training the CGN was done
104 with the default parameters that the authors also used in the paper.

105 **ImageNet counterfactuals** To produce ImageNet counterfactuals from a class conditional variable and random vector
106 we regressed a CGN with the pre-trained BigGAN backbone, as defined in Section 3.2. We investigate the quality of
107 the generative network in two ways: first by visual inspection and, secondly, measure the Inception Score and mask
108 mean μ_{mask} . Further, during training, we closely monitor the elements of the compositions to confirm the intended loss
109 behaviour. The hyperparameters for the losses and learning rates are provided by the authors. This includes the lambdas
110 defined in Equation 2 and learning rates 8E-6, 3E-5 and 1E-5 for shape, texture and background, respectively. For our
111 single GPU with 24GB memory the highest possible batch size was 5, which requires changing the number of episodes
112 and batch accumulation to 200 and 500, respectively. To this end, we regress $5 \cdot 500 \cdot 200 = 5 \cdot 10^5$ unique images
113 taking approximately 30 hours. Finally, we experiment with a modified texture mechanism where the patch grid is
114 created by filling the image randomly with the object till a degree is met.

115 **Inception Score and μ_{mask}** To assess the quality of the generated images from the ImageNet model, we calculate
116 the Inception Score⁴ as introduced by Salimans et al. (2016) for a uniform sample of non-counterfactual images. The
117 authors didn’t state in their paper how many samples they used. So, we chose to use 50,000 images with no splits as
118 Barratt and Sharma (2018) suggest that this is an appropriate amount of images given the number of classes in ImageNet.
119 The generated images from the CGN include a mask for each image of which we compute the mean pixel value μ_{mask} .

120 We calculated the inception score for our self trained CGN, the CGN using the weights provided by the authors and
121 a pretrained BigGAN model. As BigGAN does not generate masks, the μ_{mask} value was determined only for the
122 pretrained and self trained CGNs.

123 3.4.2 Claim 2

124 The following are experiments carried out to investigate the claim ‘Including generated counterfactuals in the training
125 data set improves the classifier’s out-of-domain robustness.’

⁴The following TensorFlow implementation of the Inception Score was used: <https://github.com/tsc2017/Inception-Score>

126 **Classifying MNIST datasets** Like the authors of the original paper, we trained a classifier on MNIST data and
127 compared the performance of the classifier for different compositions of the training data. We compare classifiers
128 trained on original datasets, original datasets with counterfactuals produced by the trained CGN, original dataset with
129 non-counterfactual samples generated by the CGN and only on non-counterfactual samples generated by the CGN.
130 We test on counterfactuals generated manually, so not by any CGN. We also test on non-counterfactual samples. We
131 used a classifier that has the same architecture as the one used by the authors, with the default parameters in the code
132 provided by the authors. We test how the testing accuracy changes with different datasizes and differing ratios of
133 number of counterfactuals in the training data. As an additional experiment, we produced visual explanations for the
134 classifiers trained on double coloured MNIST. These visual explanations were produced using GradCAM by Selvaraju
135 et al. (2016) applied onto the last convolutional layer of the model.

136 **ImageNet-mini classification** Similar to the authors, we measure the out-of-domain robustness for the ResNet-50
137 classifier with the ImageNet-9 background challenge dataset Xiao et al. (2020). This dataset contains two subsets
138 that hold images with randomized backgrounds of the same class and of different classes. Dubbed mixed-same and
139 mixed-rand, their difference in classification top-1 accuracy is a solid measure of class background dependence. Not
140 similar to the authors, we performed classification training on a subset of ImageNet named ImageNet-mini by Figotin
141 (2020). This dataset contains fewer images per class, significantly decreasing convergence time while only marginally
142 dropping accuracy. With the ImageNet-mini and ImageNet-9 dataset we train the ResNet-50 classifier in three ways.
143 Setting (1) contains ImageNet-mini training data only, (2) adds counterfactuals produced with the authors CGN weights,
144 and (3) adds counterfactuals produced by our weights. The amount of random counterfactuals produced for training was
145 10^5 and is suspected to be sufficient. The hyperparameters to train the classifier were searched to obtain high accuracies
146 on ImageNet-mini and are presented in 7. The search concludes on learning rate $1E-4$ and counterfactual ratio 2.0,
147 where momentum is 0.9 and weight decay is $1E-4$. Convergence point for each setting is also presented for ease of
148 verification and were all reached within 5 hours on batch size 32.

149 3.4.3 Claim 3

150 The following are experiments carried out to investigate the claim *‘Including generated counterfactuals in training data*
151 *set only marginally degrades overall accuracy’*

152 **Classifying MNIST datasets** We also tested the performance of classifiers trained on the original dataset to per-
153 formance of the classifier trained on counterfactuals (with or without original data) with non-counterfactual samples
154 as test data. This tells us whether training a classifier with counterfactuals affects performance on in-domain (non-
155 counterfactual) data.

156 **ImageNet-mini classification** Copying the experimental setting for classification training in claim 2, we test the
157 classifier on the base unmodified images of ImageNet-9 for in-domain accuracies. Additionally we inspect the top-1
158 and top-5 accuracies of our ensemble classifier on ImageNet-mini itself.

159 Besides in-domain we investigate the shape-biases of the resulting classifier ensemble with the Cue Conflict dataset
160 Gatys et al. (2015). The dataset consists of images that are generated by mixing a random texture and random object
161 in an iterative style manner. Higher bias indicates classification based on object, whereas a lower bias indicates
162 classification on texture. This helps us answer whether we can control the separate heads for shape, texture and
163 background while maintaining in-domain accuracy.

164 Since this experiment is evaluated on the same settings as claim 2, the evaluation was performed simultaneously. Not
165 increasing computation time.

166 4 Results

167 In this section we describe the results of the experiments described in the methodology. We find that our experiments
168 support claim 2 and claim 3, but we didn’t find sufficient support for claim 1 across the datasets.

169 **4.1 Claim 1**

170 In the following sections we describe results of experiments carried out to investigate the claim ‘A CGN can generate
171 high-quality counterfactual images with direct control over shape, texture, and background.’.

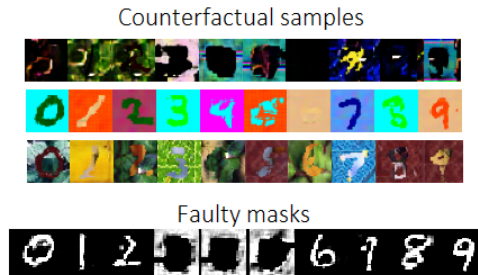


Figure 2: Samples generated by the CGN trained on the three MNIST datasets. For each of the three dataset, counterfactuals are shown. Especially colored MNIST and wildlife are of low quality. The last row shows faulty masks that were learned by the CGN.

172 **MNIST counterfactuals** By running the code for training the CGN several times, we eventually managed to train a
173 CGN that produced satisfactory samples both for the MNIST wildlife and MNIST double colored dataset. However,
174 after attempting to train the CGN at least 10 times with different initialisation weights, the results were not satisfactory
175 for the colored MNIST dataset. Some examples are shown in Figure 2. As shown in the Figure, the CGN learns faulty
176 masks for some of the digits for the colored MNIST dataset. When asked whether the authors experienced similar
177 problems training the MNIST CGN, the response was that they experienced the same, but sometimes managed to train
178 a CGN with satisfactory results. Figure 2 shows some of the faulty masks that were learned by the CGN. The results for
179 the double-colored and wildlife samples seem satisfactory, although the masks learned for the digits four and five are
180 not of high quality. The authors do not specify what percentage of trials yielded satisfactory results. We found that
181 approximately 2/3 of training trials yielded results that are satisfactory for double colored MNIST and wildlife MNIST.
182 No good training trials were found for colored MNIST. In Appendix A, more examples can be found.

183 **ImageNet counterfactuals** The quality of the compositional network is presented in Table 3 by the Inception Score
184 and mean mask value. Training the CGN ourselves obtains an IS of 115.5 compared to their trained weights 129.4,
185 indicating rather low generative capabilities relative to the BigGAN. In Appendix F we investigate the three CGN
186 mechanisms and three problems they can cause. Our IS results align well with that of the original paper 130.2. The
187 reported μ_{mask} of 0.33 indicates no mask collapse and is within the deviation shown by the authors $0.3 \pm 0.2\%$.
188 Although quality of counterfactuals are correlated with the non-counterfactual IS, we additionally investigate 16
189 randomized counterfactual images shown in Figure 7 of Appendix E. Visually the generated counterfactuals appear
190 unreal, but the quality we are after is an accurate shape, texture and background on the conjoined image. These random
191 counterfactuals appear to show these qualities.

192 **4.2 Claim 2**

193 In the following sections we describe results of experiments carried out to investigate the claim ‘Including generated
194 counterfactuals in the training data set improves the classifier’s out-of-domain robustness.’

195 **Classifying MNIST datasets** In Table 2, we show the results of training classifiers on several different datasets.
196 When trained on only the original dataset, performance on the testset is 39% for colored MNIST, but only around 10%
197 for both double-colored and wildlife. This was expected and similar to what the authors found. When trained on only
198 counterfactual data or a combination of both original data and counterfactual data, performance increases significantly.
199 While performance increases compared to performance when trained only on original data, it doesn’t reach performance
200 reported in the original paper, except for double-colored MNIST. However, performance is proportional to the quality of
201 our CGN, since especially colored MNIST and wildlife MNIST were difficult to train the CGN on (see also the results
202 in Figure 2).

203 **ImageNet-mini classification** The accuracies for ImageNet-9 background challenge are presented in Table 5. Our
 204 obtained out-of-domain robustness is highest when training on ImageNet-mini and author provided CGN weights,
 205 with a background gap of 7.7%. This cannot be compared with the authors presented value of 3.3% since a different
 206 dataset is used. Instead, we compare the difference in accuracy when including counterfactuals. This difference is
 207 -1.2% , indicating an increased out-of-domain robustness, equal to the authors reported -1.2% . Using our own weights
 208 produces lower accuracies, but the difference seems negligible.

209 **Gradient heatmaps** Figure 3 shows that when the classifier is trained on the non-counterfactual double colored
 210 dataset the positive gradients in the last convolutional layer are correlated with the color theme used in the image.
 211 However, when trained on the double colored counterfactual dataset the gradients only slightly vary when changing
 212 the color theme. We also observe that a significant part of the heatmaps from the double colored original classifier are
 213 empty, indicating that the gradient is not positive for the whole layer.

214 4.3 Claim 3

215 In the following sections we describe results of experiments carried out to investigate the claim ‘Including generated
 216 counterfactuals in training data set only marginally degrades overall accuracy’

217 **Classifying MNIST datasets** In Table 1 we show that training the classifier on counterfactuals as well as on original
 218 data only marginally decreases the accuracy on the test data when the test data consists of in-domain samples.

219 **ImageNet-mini classification** Since claim 3 is evaluated on the same training configuration of claim 2, we present
 220 the in-domain accuracies of unmodified ImageNet-9 images in Table 5. There appears no drop in accuracy when using
 221 authors weights 0.0%, which does not align with their reported drop of 1.4% on IN-9. Using are our own weights we
 222 achieve comparable results on IN-9.

223 Table 4 shows the shape bias and classifier ensemble accuracies on ImageNet-mini. When the classifier is trained with
 224 counterfactuals the texture head is able achieve a higher top-1 accuracy of $+3.4\%$ and shows an appropriately low
 225 shape bias. This shows we can individually control the three heads for classification. However, our results are not
 226 comparable with the authors since a different dataset is used. But results of shape bias are similar for our dataset.

	Colored MNIST		Double colored MNIST		Wildlife MNIST	
	Test accuracy	Train accuracy	Test accuracy	Train accuracy	Test accuracy	Train accuracy
Original	99.8	99.8	100.0	97.6	100.0	99.9
Original + CGN	99.8	94.7	99.9	94.4	99.6	99.1

Table 1: Accuracies for classifiers trained original or original and counterfactual data, and tested on test data containing in-domain images.

	Colored MNIST		Double colored MNIST		Wildlife MNIST	
	Test accuracy	Train accuracy	Test accuracy	Train accuracy	Test accuracy	Train accuracy
Original	39.0	99.8	10.1	100.0	10.6	99.4
GAN	11.7	95.4	10.0	98.7	10.8	95.1
Original + GAN	38.1	99.9	10.1	100.0	10.7	100.0
CGN	32.5	90.0	87.3	92.8	70.2	99.1
Original + CGN	53.5	91.1	87.9	94.4	63.2	97.0

Table 2: Test and training accuracy of classifiers trained on several different datasets. The testset consisted of counterfactual data. The counterfactuals used in the training data were generated by a CGN that was trained by us.

Trained on	Shape Bias	top-1 IN Acc	top-5 IN Acc
ImageNet-mini	29.1 %	65.7 %	88.2 %
IN-mini + CGN/Shape	49.6 %		
IN-mini + CGN/Text	18.0 %	69.1 %	88.1 %
IN-mini + CGN/Bg	23.1 %		

Table 4: Shape bias of the shape, texture and background classification heads. With accuracies on ImageNet-Mini.

Model	IS	μ_{mask}
CGN (theirs)	129.4	0.332
CGN (ours)	115.5	0.286
BigGAN	195.9	-

Table 3: Inception Score and mean mask μ_{mask} of CGN.

Trained on	top-1 Test Accuracies			
	IN-9	Mixed-Same	Mixed-Rand	BG-Gap
ImageNet (base)	94.7 %	85.6 %	78.1 %	7.5 %
IN-mini	91.6 %	81.8 %	73.3 %	8.5 %
IN-mini + CGN	91.6 %	81.9 %	74.2 %	7.7 %
IN-mini + our CGN	89.7 %	81.3 %	72.7 %	8.6 %

Table 5: ImageNet-9 accuracy with and without counterfactuals.

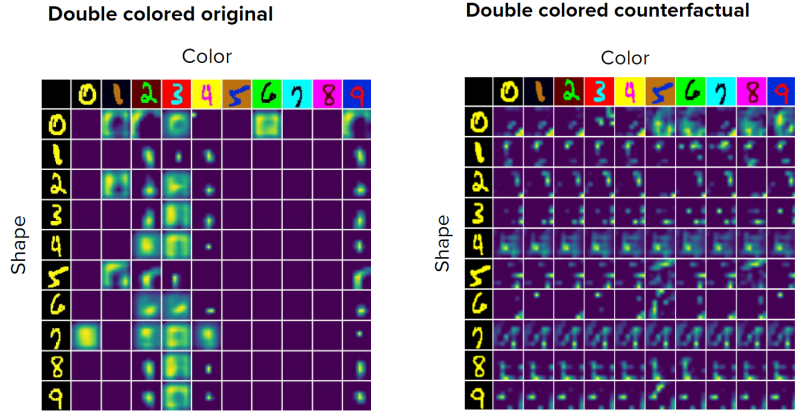


Figure 3: GradCAM heatmaps for our classifier trained on the double colored MNIST dataset. The x-position in the grid determines the color theme used for the image and the y-position determines the shape. The numbers above the x-axis correspond to the color theme used for that specific digit in the in-domain dataset.

227 5 Discussion

228 **Claim 1** The first claim, *CGNs can generate high-quality counterfactual images with direct control over shape,*
 229 *texture, and background*, was not trivial to reproduce. In particular, we found that it was difficult to train CGNs on the
 230 MNIST dataset that can generate satisfactory images. The authors of the original paper likely kept trying for longer.
 231 However, the difficulty we had with training the CGN might indicate that the architecture of the CGN can be improved
 232 to make the training process easier and get better results.

233 Our experiments with ImageNet yielded very similar and stable results. Although our quality of the generated images is
 234 slightly lower than that of the authors of the original paper. We remain within an acceptable lower inception score. The
 235 reason for this might be the training time. However, when increasing training time, failure cases such as background
 236 residue become more common, this is further discussed in Appendix F.

237 **Claim 2** The claim *Including generated counterfactuals in the training data set improves the classifier’s out-of-domain*
 238 *robustness* is supported by our experiments. For the MNIST datasets, Table ?? shows that adding counterfactual data to
 239 the training data improves accuracy on out-of-domain test data.

240 Figure 3 further supports claim 2, as we observe that the positive gradients in the final layer for the classifier trained on
 241 counterfactuals do not spuriously correlate to the color theme used in the image, in contrast to the classifier trained on
 242 the original dataset.

243 For ImageNet-mini, the decreased BG-gap shows that when adding counterfactual data to training data, the classifier’s
 244 out-of-domain robustness is increased. Although using our own weights does lead to lower scores across the test-set,
 245 which was expected due to the lower IS score of the CGN.

246 **Claim 3** The claim *‘Including generated counterfactuals in training data set only marginally degrades overall*
 247 *accuracy’* is supported by our experiments.

248 Our MNIST experiments show that adding counterfactuals barely decreases accuracy on in-domain data. The accuracies
 249 we find are lower than those of the authors. This is likely due to the fact that the quality of the CGNs that we were able
 250 to train is lower than that of the authors.

251 For ImageNet, we show no degradation in classifier accuracy on the unmodified IN-9 test set. This is not similar to the
 252 authors and is most likely caused by our choice of smaller sized ImageNet-mini subset. The CGN artificially increases
 253 dataset size and therefore mainly helps on smaller datasets. Further we have shown similar biases across the classifier
 254 ensemble indicating full control off the classifier decision making. With the texture classification head we report an
 255 increase in top-1 accuracy on in-domain ImageNet-mini, but is also most likely caused by the artificial increase in
 256 dataset size.

257 **Bibliography**

258 Barratt, S. and Sharma, R. (2018). A note on the inception score.

259 Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., and Vedaldi, A. (2013). Describing textures in the wild. *CoRR*,
260 abs/1311.3618.

261 dos Santos Tanaka, F. H. K. and Aranha, C. (2019). Data augmentation using gans.

262 Figotin, I. (2020). Imagenet-1k-mini dataset. <https://www.kaggle.com/ifigotin/imagenetmini-1000>.

263 Gatys, L. A., Ecker, A. S., and Bethge, M. (2015). A neural algorithm of artistic style.

264 Goyal, M., Rajpura, P., Bojinov, H., and Hegde, R. (2018). Dataset augmentation with synthetic images improves
265 semantic segmentation. *Computer Vision, Pattern Recognition, Image Processing, and Graphics*, page 348–359.

266 Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. (2016). Improved techniques for
267 training gans.

268 Sauer, A. and Geiger, A. (2021). Counterfactual generative networks. In *International Conference on Learning
269 Representations*.

270 Selvaraju, R. R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., and Batra, D. (2016). Grad-cam: Why did you say
271 that? visual explanations from deep networks via gradient-based localization. *CoRR*, abs/1610.02391.

272 Wong, S. C., Gatt, A., Stamatescu, V., and McDonnell, M. D. (2016). Understanding data augmentation for classification:
273 When to warp? In *2016 International Conference on Digital Image Computing: Techniques and Applications
274 (DICTA)*, pages 1–6.

275 Xiao, K., Engstrom, L., Ilyas, A., and Madry, A. (2020). Noise or signal: The role of image backgrounds in object
276 recognition.

278 A MNIST counterfactuals generated by our CGNs



Figure 4: **MNIST counterfactuals.** From left to right: colored, double-colored and wildlife MNIST. Images are generated by the CGNs which were trained by us.

279 B MNIST classification with weights of authors original paper

	colored MNIST		double-colored MNIST		wildlife MNIST	
	Train Acc	Test Acc	Train Acc	Test Acc	Train Acc	Test Acc
Original	99.7 %	38.9 %	100.0 %	10.1 %	99.3 %	10.6 %
GAN	99.9 %	25.83 %	100.0 %	9.9 %	100.0 %	10.8 %
Original + GAN	99.8 %	40.3 %	100.0 %	10.0 %	100.0 %	10.7 %
CGN	99.3 %	92.8 %	96.7 %	90.2 %	98.5 %	84.4 %
Original + CGN	99.2 %	96.4 %	97.2 %	87.2 %	98.0 %	76.24 %

Table 6: Test and training accuracy of classifiers trained on several different datasets. The testset consisted of counterfactual data. The counterfactuals used in the trainingdata were generated by a CGN with weights provided by the authors of the original paper.

280 **C MNIST ablation study**

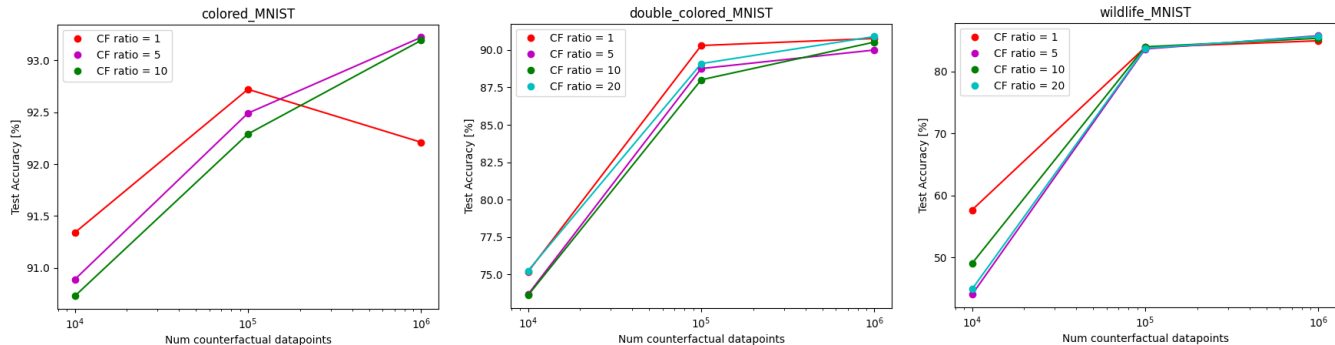


Figure 5: Test accuracy for classifiers trained on original data and counterfactual data with different CF ratios. The CF ratio indicates how many counterfactuals we generate per sampled noise. For colored MNIST, the maximum CF ratio is ten as there are only ten possible colors per shape. The counterfactuals in the trainingdata were generated by a CGN we trained ourselves.

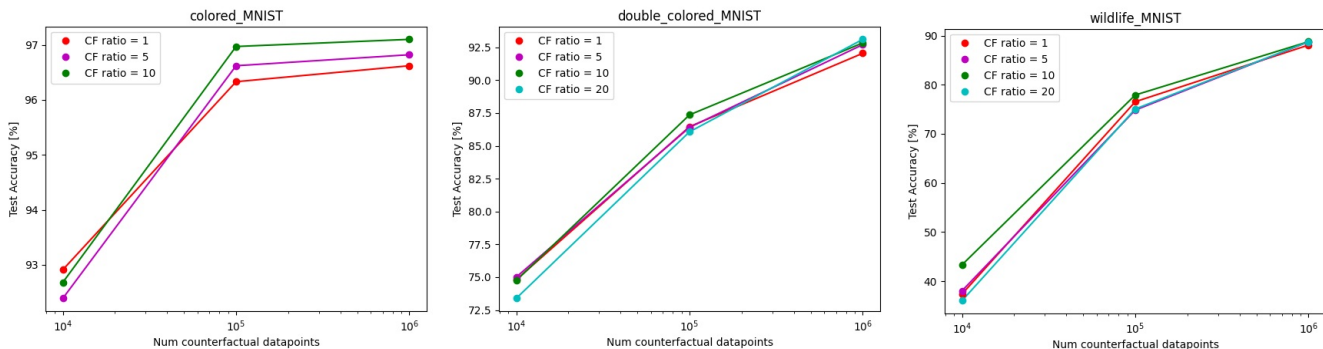


Figure 6: Test accuracy for classifiers trained on original data and counterfactual data with different CF ratios. The CF ratio indicates how many counterfactuals we generate per sampled noise. For colored MNIST, the maximum CF ratio is ten as there are only ten possible colors per shape. The counterfactuals in the trainingdata were generated by a CGN with weights provided by the authors.

281 **D ImageNet hyperparameter search**

Trained on	Lr	CF-ratio	Top-1 Test Accuracies				
			Epoch	IN-9	Mixed-Same	Mixed-Rand	BG-Gap
ImageNet (base)		0.0	0	94.7 %	85.6 %	78.1 %	7.5 %
IN-mini	1E-4	0.0	16	91.6 %	81.8%	73.3 %	8.5 %
	1E-3	0.0	22	90.7 %	79.6 %	69.6 %	10.0 %
IN-mini + CGN	1E-4	1.0	11	90.9 %	81.8 %	73.5 %	8.3 %
	1E-4	2.0	20	91.6 %	81.9 %	74.2 %	7.7 %
	1E-3	1.0	18	88.7 %	79.9 %	71.4 %	8.5 %
IN-mini + our CGN	1E-3	2.0	28	88.7 %	78.9 %	70.6 %	8.3 %
	1E-4	1.0	17	89.7 %	81.3 %	72.7 %	8.6 %
	1E-4	2.0	16	91.4 %	81.7 %	72.3 %	9.4 %

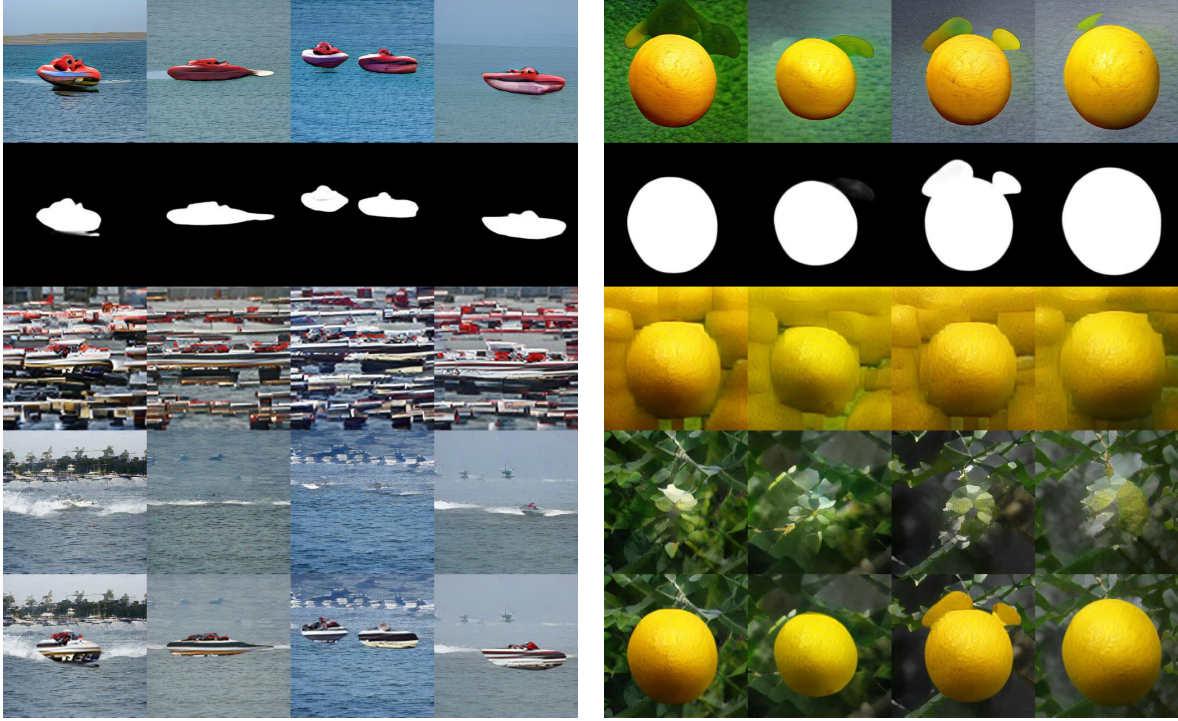
Table 7: **ImageNet-9 classification hyperparameters.** Investigated hyperparameters for the pre-trained ResNet-50 model.



Row	Column	Shape	Texture	Background
1	1	vase	golden retriever	black and gold garden spider
1	2	wine bottle	king penguin	beaver
1	3	hartebeest	red wolf	drake
1	4	wardrobe	bookshop	valley
2	1	mushroom	agaric	breakwater
2	2	toaster	jack-o'-lantern	jackfruit
2	3	tub	ambulance	valley
2	4	drake	megalith	paddle

Row	Column	Shape	Texture	Background
3	1	tripod	fire engine	beaver
3	2	sweatshirt	head cabbage	cliff
3	3	drake	tennis ball	breakwater
3	4	ostrich	golden retriever	geyser
4	1	tripod	plastic bag	rock crab
4	2	ostrich	monarch butterfly	grey whale
4	3	orange	analog clock	viaduct
4	4	red wine	golden retriever	ostrich

Figure 7: **Counterfactual images.** Here, counterfactual images generated by the CGN, that we trained, are displayed. In the table below the image, labels given to each independent mechanism can be found.



(a) IM outputs for 'boat'. From top to bottom: \tilde{m} , m , f , b , x_{gen}

(b) IM outputs for 'lemon.' From top to bottom: \tilde{m} , m , f , b , x_{gen}

Figure 8: **IM outputs:** Output of the each mechanism in Fig 1 is displayed. From top to bottom; pre-trained BigGAN, f_{shape} , $f_{texture}$, $f_{background}$ output and finally the composition mechanism's output (x_{gen}) is displayed.

283 **F Mechanism and failures of ImageNet CGN**

284 In this section, we discuss the three supposedly independent components of the CGN and their failure cases.

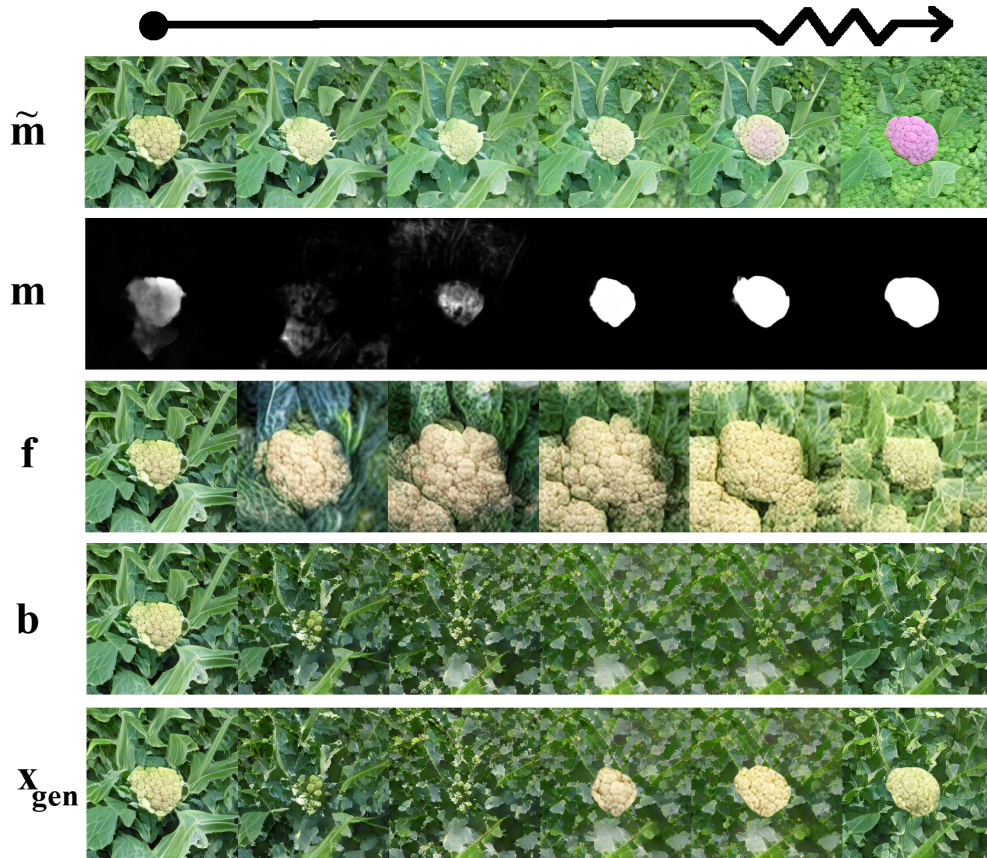


Figure 9: **IM Outputs for cauliflower** Learnt pre-masks \tilde{m} , masks m , foregrounds f , and backgrounds b . The arrow indicates the beginning of training till the jumped ending.

285 The figure above indicates visually that the individual losses cause intended changes in the BigGAN for each mechanism.
 286 Pre-mask \tilde{m} accentuates the location of the object, foreground f shows near complete texture mapping, some background
 287 artifacts remain in the texture which occurs often across classes. Mask b always converges well with no remaining
 288 object across the classes, although object artifacts do occur when training the network longer.

289 Failure cases of the CGN are displayed under three main categories; texture-background entanglement, background
 290 residue and reduced realism. All those images are generated using the CGN we trained.

291 We conjecture that most texture-background entanglement occurs when the mask of the object is either difficult to learn,
 292 is relatively small $\mu_{mask} < 0.1$ with $\tau = 0.1$, or the patch size does not match the object well. All of these cases cause
 293 the object mask to include more background for lower loss, resulting in texture-background entanglement.

294 Background residue occurs when background loss has converged during training. The background BigGAN then further
 295 improves by decreasing the reconstruction loss. In increasing the background details, that are outside the mask, the
 296 inside of the mask is allowed to change in all ways as long as no object is detected. This results in undetected artifacts
 297 that are increasingly present in longer trained models.

298 With a perfect object mask, texture and background identical realism to the BigGAN is expected. The visual results and
 299 Inception Score confirm that this is not the case. Reduced realism is prevalent amongst all generated images and are
 300 conjectured to be caused by several issues. Which are an entangled problem between the mechanisms.

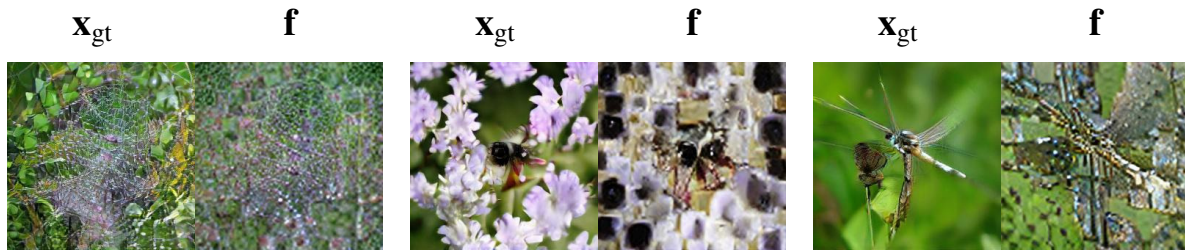


Figure 10: **Texture-Background Entanglement.** Texture maps contain traces from the background. From left to right, the images are; spiderweb, bee and dragonfly.

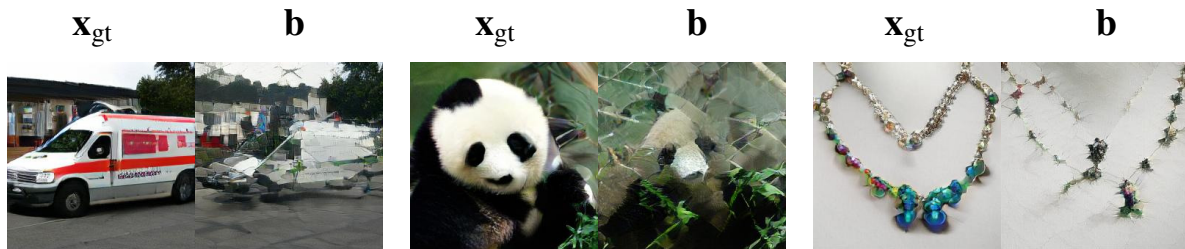


Figure 11: **Background Residue.** Regions, where the objects are located, are not fully in-painted. Background still contains some artifacts from the object.

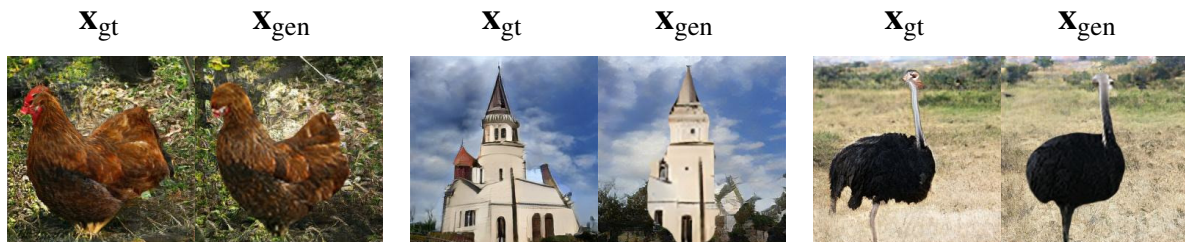


Figure 12: **Reduced realism.** Generated images mostly do not look realistic. Authors stated that this is due to the constraints enforced and the analytically defined composition mechanism.