# PROREGEN: PROGRESSIVE RESIDUAL GENERATION UNDER ATTRIBUTE CORRELATIONS

# **Anonymous authors**

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027 028 029

031

032

033

034

037

040

041

042

043

044

046

047 048

051

052

Paper under double-blind review

# **ABSTRACT**

Attribute correlations in the training data will compromise the ability of a deep generative model (DGM) to synthesize images with under-represented attribute combinations (i.e., minority samples). Existing approaches mitigate this by data re-sampling to remove attribute correlations seen by the DGM, using a classifier to provide pseudo-supervision on generated counterfactual samples, or incorporating inductive bias to explicitly decompose the generation into independent submechanisms. We present ProReGen, a progressive residual generation approach inspired by the classical Robinson's transformation, to partial out from an image attribute  $x_2$  its component  $m(x_1)$  that is predictable by other image attributes  $x_1$ , and the residual  $\gamma = \mathbf{x}_2 - m(\mathbf{x}_1)$  that is not. This simplifies the problem of learning a DGM  $g(\mathbf{x}_1, \mathbf{x}_2)$  conditioned on correlated inputs, to learning  $\tilde{g}(\mathbf{x}_1, \gamma)$ conditioned on orthogonal inputs. It further allows us to progressively learn  $\tilde{g}$  by first shifting the burden to abundant majority samples to learn  $\tilde{g}(\mathbf{x}_1, \gamma = 0)$ , and then expanding it with additional layers  $g_{res}$  to resolve its difference to  $\tilde{g}(\mathbf{x}_1, \gamma)$  using residual attribute  $\gamma$  on limited minority samples. On three benchmark datasets with varying strengths of attribute correlations, we demonstrate that ProReGen with input orthogonalization and progressive residual learning—improved the correctness of minority generations compared to existing strategies.

# 1 Introduction

Attribute correlations are not uncommon in observed image datasets. Some may be a natural manifestation of underlying causal relations, *e.g.*, the object in an image determining the background (Sagawa et al., 2019). Some may reflect bias in data curation, *e.g.*, collecting patient data from those who already received treatment (Wang et al., 2017). Regardless of the mechanisms, attribute correlations can induce unintended consequences in deep neural networks (DNN) training.

In the context of discriminative (e.g., classification) DNNs, this phenomenon has been widely discussed, often under the concept of *spurious correlations* or *short-cut learning* (Ye et al., 2024). In the context of deep generative models (DGM), such discussion is comparatively less structured and scatters across a variety of topics. On one hand, the importance for a DGM to properly synthesize under-represented image examples—those with image attributes that do not comply with the observed correlation—are appreciated in many domains, e.g., for explaining whether a DNN classifier has captured correlated features for decision making (Rodríguez et al., 2021), or for augmenting training data to mitigate correlations (Kim et al., 2021). On the other hand, several evaluation studies (Träuble et al., 2021; Bose et al., 2022) have shown that naively-trained DGMs would capture latent attribute correlations from training data (Träuble et al., 2021) and even reveal the associated causal directions (Bose et al., 2022). How does this impact the synthesis of under-represented samples, and to what extent could it be mitigated? Answers to these questions remain largely open.

Consider two sets of image attributes  $\mathbf{x}_1$  and  $\mathbf{x}_2$  (both can be multi-dimensional) that exhibit a correlation in observed data. Consider the goal of learning a DGM g conditioned on these attributes to generate an image  $\mathbf{y}$  as  $\mathbf{y} = g(\mathbf{x}_1, \mathbf{x}_2)$ . For the function g to generate with different combinations of  $\mathbf{x}_1$  and  $\mathbf{x}_2$  values, it is important for g to correctly model the mechanisms,  $g_1$  and  $g_2$ , through which  $\mathbf{x}_1$  and  $\mathbf{x}_2$  influences  $\mathbf{y}$  separately. Unfortunately, due to the observed  $\mathbf{x}_1$ - $\mathbf{x}_2$  correlation,  $g_1$  and  $g_2$  can only be separately observed in the small number of samples where such correlation does not hold. We stress this as a fundamental challenge for learning a DGM under attribute correlations.

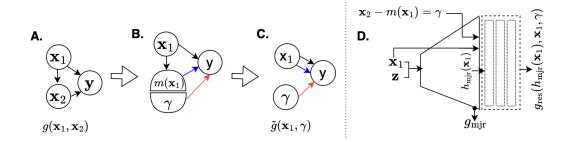


Figure 1: A-C: Overview of the Robinson's partialling-out approach in motivating learning  $g(\mathbf{x}_1, \mathbf{x}_2)$  as  $\tilde{g}(\mathbf{x}_1, \gamma)$ . D: Illustration of ProReGen to realize function  $\tilde{g}(\mathbf{x}_1, \gamma)$  by progressively learning  $g_{\text{mjr}}(\mathbf{x}_1) \coloneqq \tilde{g}(\mathbf{x}_1, \gamma = 0)$  from majority samples, followed by learning  $g_{\text{res}}(\mathbf{h}_{\text{mjr}}(\mathbf{x}_1), \mathbf{x}_1, \gamma)$  on minority samples to resolve the difference to  $\tilde{g}(\mathbf{x}_1, \gamma)$  with  $\gamma$ .

Existing strategies to address this challenge are limited. Re-sampling is a simple approach to balance training samples (Monteiro et al., 2022), essentially up-weighting under-represented samples where  $g_1$  and  $g_2$  can be separately observed. Alternatively, inductive bias has been introduced to explicitly decompose g into independent mechanisms  $g_1$  and  $g_2$ , which requires prior knowledge about how  $\mathbf{x}_1$  and  $\mathbf{x}_2$  may influence  $\mathbf{y}$  differently (e.g., object shape vs. texture vs. background) (Sauer & Geiger, 2020). Finally, to go beyond the limits of factual under-represented samples, pseudo-supervision on generated counterfactual images has been presented, typically realized by using a classifier to recognize feature attributes in the generated images (Kocaoglu et al., 2017; Ribeiro et al., 2023; He et al., 2019). However, since the classifier is trained under the same attribute correlations, its ability to correctly recognize these attributes is likely compromised – how does this impact its validity to supervise the generation of under-represented counterfactual images has not been well understood.

In this paper, we take a fundamentally different perspective to address the challenge of modeling  $g(\mathbf{x}_1, \mathbf{x}_2)$  under attribute correlations, inspired by the classical Robinson's partialling-out transformation (Robinson, 1988). While details of this concept will be introduced in Section 3, Fig. 1 illustrates its core concept in the context of modeling  $g(\mathbf{x}_1, \mathbf{x}_2)$ . Consider the causal graph in Fig. 1A with a causal direction assumed between  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . Instead of attempting to model the independent causal mechanisms  $g_1$  and  $g_2$  in the presence of such correlations, we decompose  $\mathbf{x}_2$  into  $E[\mathbf{x}_2|\mathbf{x}_1]$  that can be predicted by  $\mathbf{x}_1$ , along with a residual  $\gamma$  that cannot, i.e.,  $\mathbf{x}_2 = E[\mathbf{x}_2|\mathbf{x}_1] + \gamma$  (Fig. 1B). With this, instead of modeling  $\mathbf{y} = g(\mathbf{x}_1, \mathbf{x}_2)$  as a composition of  $g_1$  and  $g_2$  as in Fig. 1A, we model it as  $\mathbf{y} = \tilde{g}(\mathbf{x}_1, \gamma)$  as in Fig. 1C: the effect of  $\mathbf{x}_1$  on  $\mathbf{y}$  now absorbs the effect from  $E[\mathbf{x}_2|\mathbf{x}_1]$ , the component of  $\mathbf{x}_2$  that can be predicted from  $\mathbf{x}_1$  – we referred to this as correlated effect; as such, the effect of  $\gamma$  on  $\mathbf{y}$ —the part of  $\mathbf{x}_2$ 's influence on  $\mathbf{y}$  that cannot be explained by  $\mathbf{x}_1$ —is partialled out: we refer to this as residual effect. In this new causal graph (Fig. 1C), instead of attempting to recover two independent mechanisms from correlated inputs  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , we transform the problem into learning independent correlated and residual effects from two independent inputs  $\mathbf{x}_1$  and  $\gamma$ .

To design a DGM based on Fig. 1C, we note that  $\tilde{g}(\mathbf{x}_1,\gamma)$  at  $\gamma=0$ , i.e.,  $E[\mathbf{y}|\mathbf{x}_1,\gamma=0]$ , can be estimated from abundant samples whose attribute values meet the observed correlations (henceforth referred to as majority samples).  $\tilde{g}(\mathbf{x}_1,\gamma)$  at  $\gamma\neq0$ , on the other hand, can only be estimated from a limited number of samples where such correlation does not hold (henceforth referred to as minority samples). To shift the primary burden of learning  $\tilde{g}$  to majority samples, therefore, we further decompose the learning of  $\tilde{g}(\mathbf{x}_1,\gamma)$  into the learning of  $\tilde{g}(\mathbf{x}_1,0)$  at  $\gamma=0$  using majority samples, and use minority samples to make up the difference between  $\tilde{g}(\mathbf{x}_1,0)$  to  $\tilde{g}(\mathbf{x}_1,\gamma)$ . This results in our progressive residual effect generator (ProReGen) that is progressively learned in two stages as outlined in Fig. 1D. ProReGen as described is expected to have two major benefits. First, the orthogonalization of the inputs  $\mathbf{x}_1$  and  $\gamma$  helps separate the learning of their independent effects on  $\mathbf{y}$ . Second, with the progressive expansion from  $g_{\text{mjr}}$  to  $g_{\text{res}}$ , the challenge of learning to generate under attribute correlations is reduced to learning the residual between  $\tilde{g}(\mathbf{x}_1,\gamma=0)$  and  $\tilde{g}(\mathbf{x}_1,\gamma)$ . We instantiate the concept of ProReGan in conditional-VAEs and -GANs and, on three benchmark datasets with varying strengths of attribute correlations, we experimentally demonstrate the improved per-

formance of ProReGen in generating *correct* minority samples compared to naive DGMs or those strengthened with re-weighting of factual samples or pseudo-supervision on generated samples.

#### 2 RELATED WORKS

Several domains find uses for synthesizing minority images. For instance, to explain if a classifier has captured attribute correlations from its training data, one can test if it is able to guide a pre-trained DGM to generate minority images (Rodríguez et al., 2021; Jeanneret et al., 2022). In addition to explaining, synthesized minority images can also be used for augmenting and removing attribute correlations in training data (Goel et al., 2020; Kim et al., 2021). While naively-trained DGMs have been used for such syntheses (Rodríguez et al., 2021; Jeanneret et al., 2022), there is an increasing attention on the impact of attribute correlations on DGM training and potential mitigation solutions.

**Re-weighting of factual samples:** The concept of *simulated intervention* was presented in Monteiro et al. (2022) to re-sample data according to the marginal distribution of image attributes, to effectively remove attribute correlations seen by the DGM. This is essentially similar to re-weighting, where minority samples are up-weighted in their contribution to the training signals. This approach is ultimately affected by the quantity and diversity of factual minority samples.

**Pseudo-labeling of counterfactual generations:** Instead of relying on factual minority samples, an alternative is to provide some *pseudo-supervision* to encourage the DGM to generate *counterfactual* images with intended feature attributes. This is often achieved by leveraging another DNN classifier, often trained from the same data as the DGM, to provide supervisory signals by recognizing the attributes of generated images (Kocaoglu et al., 2017; Ribeiro et al., 2023; He et al., 2019). However, because the DNN classifier is also subject to attribute correlations in the training data, their reliability in correctly recognizing these attributes is questionable. How does this affect the supervisory signal it provides to the DGM's counterfactual generations has not been systematically investigated.

**Inductive bias to decompose generation mechanisms:** An entirely different approach is to incorporate inductive bias about the mechanisms under which different attributes contribute to generated images (Sauer & Geiger, 2020; Park et al., 2020) In (Sauer & Geiger, 2020), for instance, the image generation process is decomposed into independent shape, texture, and background mechanisms. While this approach tends to be highly successful when assumptions of the underlying generation mechanisms are met, it does require prior knowledge for the design of such inductive bias.

ProReGen represents a completely different approach to these existing works. Inspired by the classical Robinson's partialling out approach, ProReGen tackles the challenge of attribute correlations at its core by first recasting the DGM from being conditioned on correlated attributes to orthogonal attributes. By a progressive expansion design, it further leverages majority training samples to reduce the problem of learning minority generation to learning its residual to majority generation. The latter concept is related to existing works in data augmentation, where image translation models are used to transform a factual majority sample to a minority counterfactual. Examples include the use of swapping-autoencoder (Kim et al., 2021), CycleGAN (Goel et al., 2020), and few-shot adaptation of GANs (An et al., 2022). Because these models are intended only for translating factual samples but not general image generations, they are out of the scope of our consideration.

# 3 Preliminary: Robinson's Partialling-Out Transformation

Consider a partial linear equation  $E[y|x_1, x_2] = \theta(x_1) + \beta x_2$ , where  $x_2 = m(x_1) + u$ , Robinson's transformation in (Robinson, 1988) decomposes the original equation into:

$$E[y|x_1, x_2] = \theta(x_1) + \beta * (m(x_1) + u) = E[y|x_1] + \beta * (x_2 - m(x_1))$$
(1)

where  $E[y|x_1] = \theta(x_1) + \beta * m(x_1)$ . Effectively, instead of describing the separate effect of  $x_1$  and  $x_2$  on y as dictated the original equation, we decompose their effect into two orthogonal components: 1)  $E[y|x_1]$  that describes the combined effect of  $x_1$  and  $m(x_1)$  on y—the latter absorbing the effect of  $x_2$  on  $x_2$  that can be predicted by  $x_2$ ; and 2) the *residual* effect of  $x_2$  or  $x_2$  that can be predicted by  $x_2$ ; and 2) the residual effect of  $x_2$  that cannot be predicted from  $x_2$ .

In the original paper (Robinson, 1988), the goal of this decomposition is to estimate parameter  $\beta$ , which arrives at the classical *residual-on-residual* least-square fitting of  $\beta$  via:  $E[y|x_1, x_2]$  –

 $E[y|x_1] = \beta * (x_2 - E[x_2|x_1])$ . The resulting estimator for  $\beta$  can further be shown to meet the *Neyman's orthogonality condition* such that it is insensitive to perturbations in the estimator for  $E[y|x_1]$  or  $E[x_2|x_1]$  (Robinson, 1988; Chernozhukov et al., 2018). This decomposition has also served as the basis for the R-learner (R for Robinson) that extended the approach to estimating the function  $\beta(x_1)$  instead of the low-dimensional parameter  $\beta$  (Nie & Wager, 2021).

In contrast, we use this transformation as the foundation inspiring the design of an image-generating DGM, where y is high-dimensional and its relation with  $x_1$  and  $x_2$  is highly nonlinear. Note that the causal direction  $x_1 \to x_2$  depends on the target parameter  $\beta$  in the original formulation: in our setting where y as a function of g remains the primary interest, the assumed causal direction will influence the decomposition, but not the general modeling strategy. The difficulty of the residual generation task however may change with this direction, which we will empirically examine in Section 5.4. Below we continue our discussion with  $x_1 \to x_2$  without loss of generality.

# 4 METHODOLOGY

# 4.1 PROREGEN: PROGRESSIVE RESIDUAL GENERATION

Consider a conditional-DGM  $\mathbf{y} = g(\mathbf{z}, \mathbf{x}_1, \mathbf{x}_2)$  underlying observed image data, where a correlation exists between image attributes  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , and  $\mathbf{z}$  represents latent variables not included in  $\mathbf{x}_1$  and  $\mathbf{x}_2$ . Built on Robinson's partialling-out approach, as illustrated in Fig. 1A-C, we first cast the problem of learning  $\mathbf{y} = g(\mathbf{z}, \mathbf{x}_1, \mathbf{x}_2)$  with correlated inputs, to learning  $\mathbf{y} = \tilde{g}(\mathbf{z}, \mathbf{x}_1, \gamma)$  with independent inputs where  $\gamma = \mathbf{x}_2 - m(\mathbf{x}_1)$  represents the residual in  $\mathbf{x}_2$  that cannot be predicted by  $\mathbf{x}_1$ . Because  $\tilde{g}(\mathbf{z}, \mathbf{x}_1, \gamma = 0)$  is described by majority samples vs.  $\tilde{g}(\mathbf{z}, \mathbf{x}_1, \gamma)$  at  $\gamma \neq 0$  by a small number of minority samples, we further design a progressive learning strategy to shift the burden of learning  $\tilde{g}$  mostly to the learning of  $\tilde{g}(\mathbf{z}, \mathbf{x}_1, \gamma = 0)$  using majority samples, and using minority samples to only resolve its difference to  $\tilde{g}(\mathbf{z}, \mathbf{x}_1, \gamma)$  with  $\gamma$ . This gives the foundation of ProReGen:

$$\mathbf{y} = \tilde{g}(\mathbf{z}, \mathbf{x}_1, \gamma) \approx g_{\text{res}}(\mathbf{h}_{\text{mjr}}(\mathbf{x}_1), \mathbf{x}_1, \gamma), \text{ where } \gamma = \mathbf{x}_2 - m(\mathbf{x}_1)$$
 (2)

where  $\mathbf{h}_{\text{mjr}}$  is the feature map of  $g_{\text{mjr}} \coloneqq \tilde{g}(\mathbf{z}, \mathbf{x}_1, \gamma = 0)$  before the final activation layer. Equation (2) includes three main components progressively learned in two stages, as illustrated in Fig. 1D:

- In stage-I, from a large number of majority samples, we learn an attribute predict function  $\mathbf{x}_2 = m(\mathbf{x}_1)$  to approximate  $E[\mathbf{x}_2|\mathbf{x}_1]$ , and a generative model  $g_{\text{mjr}}(\mathbf{x}_1)$  to approximate  $\tilde{g}(\mathbf{x}_1, \gamma = 0)$ , the latter effectively describing the generation when  $\mathbf{x}_2 = m(\mathbf{x}_1)$ , i.e., for a majority sample.
- In stage-II, using available minority samples, we expand the generator with additional layers,  $g_{\text{res}}$ , to resolve the residual between  $\tilde{g}(\mathbf{x}_1, \gamma)$  and  $g_{\text{mjr}}(\mathbf{x}_1)$  with the residual  $\gamma$  partialled out from  $\mathbf{x}_2$ . Effectively, we approximate  $\tilde{g}(\mathbf{x}_1, \gamma)$  by  $g_{\text{res}}(\mathbf{h}_{\text{mjr}}(\mathbf{x}_1), \mathbf{x}_1, \gamma)$ , with  $\mathbf{h}_{\text{mjr}}$  defined above.

This learning process takes form of a progressive-DGM, where a backbone  $g_{\rm mjr}$  is first learned on majority samples and then expanded on minority samples. This concept is agnostic to the type of DGMs: below, we describe its instantiations on conditional-VAEs (c-VAEs) and -GANs (c-GANs).

#### 4.2 PROREGEN-VAE

Stage I: Learning a c-VAE that captures attribute correlations: To learn  $g_{\rm mjr}({\bf z},{\bf x}_1)$  as a c-VAE, we define a decoder network  $G_{\theta_{\rm mjr}}({\bf z},{\bf x}_1)$  that parameterizes the likelihood  $p_{\theta_{\rm mjr}}({\bf y}\mid{\bf z},{\bf x}_1)$ , and its corresponding encoder network  $E_{\phi_{\rm mjr}}({\bf y},{\bf x}_1)$  that parameterizes the approximate posterior  $q_{\phi_{\rm mjr}}({\bf z}\mid{\bf y},{\bf x}_1)$ , both conditioned on attribute labels  ${\bf x}_1$ . They are trained on majority samples by maximizing the standard ELBO loss:

$$\max_{\theta_{\rm mjr},\phi_{\rm mjr}} \left\{ \underbrace{\mathbb{E}_{\mathbf{z} \sim E_{\phi_{\rm mjr}}(\mathbf{y}_{\rm mjr},\mathbf{x}_1)} \left[ - \|\mathbf{y}_{\rm mjr} - G_{\theta_{\rm mjr}}(\mathbf{z},\mathbf{x}_1)\|_2^2 \right]}_{\text{Reconstruction Loss}} - \beta \underbrace{D_{\rm KL} \left( E_{\phi_{\rm mjr}}(\mathbf{y}_{\rm mjr},\mathbf{x}_1) \| p(\mathbf{z}) \right) \right\}}_{\text{KL Divergence}}$$
(3)

where  $p(\mathbf{z})$  is defined as the standard isotropic Gaussian prior  $\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ , and hyperparameter  $\beta > 0$  adjusts the KL-regularization strength. While only conditioned on  $\mathbf{x}_1$ ,  $G_{\theta_{\text{mjr}}}$  is expected to absorb the effect of  $m(\mathbf{x}_1)$ , the portion of  $\mathbf{x}_2$  predictable by  $\mathbf{x}_1$ , due to their natural correlation as observed. At the same time, a function  $\mathbf{x}_2 = \hat{m}(\mathbf{x}_1)$  is estimated on the attribute values from majority samples.

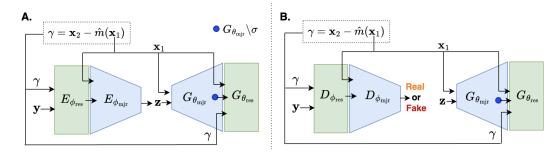


Figure 2: Illustration of ProReGen-VAE (left) and ProReGen-GAN (right)

Stage II: Learning the residual by expanding the c-VAE: To leverage the learned  $g_{\rm mjr}({\bf z},{\bf x}_1)$ , we now expand the decoder  $G_{\theta_{\rm mjr}}$  with several additional layers at the last layer of  $G_{\theta_{\rm mjr}}$  before the final activation layer  $\sigma$ , which we denote as  $G_{\theta_{\rm mjr}}\backslash\sigma$ . We denote these expanded portion of the decoder as  $G_{\theta_{\rm res}}(G_{\theta_{\rm mjr}}\backslash\sigma,{\bf x}_1,\gamma)$ , conditioned on  $\gamma={\bf x}_2-\hat{m}({\bf x}_1)$  that represents the residual in  ${\bf x}_2$  that cannot be predicted by  ${\bf x}_1$ . As illustrated in Fig. 2A, we expand the encoder network with additional layers  $E_{\phi_{\rm res}}$  in a mirror of the expanded decoder network as  $E_{\phi_{\rm res}}({\bf y}_{\rm mnr},\gamma)$ , to produce output that will serve as the input to the first-stage encoder. The expanded networks are trained on minority samples by maximizing the ELBO loss below, where we keep the stage-I weights  $\theta_{\rm mjr}$  and  $\phi_{\rm mjr}$  frozen:

$$\max_{\theta_{\text{res}},\phi_{\text{res}}} \{ \mathbb{E}_{\mathbf{z} \sim E_{\phi_{\text{mjr}}}(E_{\phi_{\text{res}}}(\mathbf{y}_{\text{mnr}},\gamma),\mathbf{x}_1)} [\|\mathbf{y}_{\text{mnr}} - \hat{\mathbf{y}}_{\text{mnr}}\|_2^2] + \beta \underbrace{D_{\text{KL}} (E_{\phi_{\text{mjr}}}(E_{\phi_{\text{res}}}(\mathbf{y}_{\text{mnr}},\gamma),\mathbf{x}_1) \| p(\mathbf{z}))}_{\text{KL Divergence}} \}$$
(4)

where 
$$\hat{\mathbf{y}}_{mnr} = G_{\theta_{res}}(G_{\theta_{mjr}\setminus\sigma}(\mathbf{z},\mathbf{x}_1),\mathbf{x}_1,\gamma), \quad \mathbf{z} \sim E_{\phi_{mjr}}(E_{\phi_{res}}(\mathbf{y}_{mnr},\gamma),\mathbf{x}_1)$$
 (5)

While Equation (5) represents a general formulation for residual generation, on simpler datasets, an additive residual can be considered such that  $\hat{\mathbf{y}}_{\text{mnr}} = G_{\theta_{\text{mjr}}}(\mathbf{z}, \mathbf{x}_1) + G_{\theta_{\text{res}}}(G_{\theta_{\text{mjr}}} \setminus_h(\mathbf{z}, \mathbf{x}_1)), \gamma)$ . With Equation (5),  $G_{\theta_{\text{res}}}$  leverages the feature map of the stage-I generator  $G_{\theta_{\text{mjr}}}$  and the residual  $\gamma$  to resolve the difference between majority and minority generations. Intuitively, because  $G_{\theta_{\text{mjr}}}$  and  $E_{\phi_{\text{mjr}}}$  are trained in stage-I to generate and encode from majority samples, the expanded  $G_{\theta_{\text{res}}}$  will be encouraged to learn to modify a majority-image feature map to include features corresponding to the residual partialled-out from  $\mathbf{x}_2$ , while the expanded  $E_{\theta_{\text{res}}}$  will be encouraged to alter such features to generate an output feature map acceptable to  $E_{\phi_{\text{mjr}}}$  (i.e., compliant with feature map seen by  $E_{\phi_{\text{mjr}}}$  in stage-I). With this, we shift the burden of learning the c-VAE mainly to majority samples, and allow the use of limited minority samples for learning the necessary reisdual changes only.

#### 4.3 PROREGEN-GAN

Stage-I: Learning a c-GAN that captures attribute correlations: To learn  $g_{mjr}(\mathbf{z}, \mathbf{x}_1)$  as a c-GAN, we define a generator  $G_{\theta_{mjr}}(\mathbf{z}, \mathbf{x}_1)$  and its discriminator  $D_{\phi_{mjr}}(\mathbf{y}, \mathbf{x}_1)$ , both conditioned on attribute labels  $\mathbf{x}_1$ . They are trained on majority samples using standard adversarial loss:

$$\min_{\theta_{\text{mjr}}} \max_{\phi_{\text{mjr}}} \left\{ \mathbb{E}_{\mathbf{y}_{\text{mjr}}, \mathbf{x}_{1} \sim p_{\text{data}}} \left[ \log D_{\phi_{\text{mjr}}}(\mathbf{y}_{\text{mjr}}, \mathbf{x}_{1}) \right] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z}), \mathbf{x}_{1} \sim p_{\text{data}}} \left[ \log \left( 1 - D_{\phi_{\text{mjr}}}(G_{\theta_{\text{mjr}}}(\mathbf{z}, \mathbf{x}_{1}), \mathbf{x}_{1}) \right) \right] \right\}$$
(6)

where  $D_{\phi_{\rm mjr}}$  is trained to maximize the probability of correctly distinguishing the real  $\mathbf{y}_{\rm mjr}$  vs. samples generated by  $G_{\theta_{\rm mjr}}(\mathbf{z}, \mathbf{x}_1)$ ; while the generator  $G_{\theta_{\rm mjr}}$  is trained to fool the discriminator. Similar to ProReGen-VAE, while only conditioned on  $\mathbf{x}_1$ , the generator  $G_{\theta_{\rm mjr}}$  is expected to absorb the effect of  $m(\mathbf{x}_1)$  and generate samples representative of the majority samples in the observed data. In the meantime, a function  $\mathbf{x}_2 = \hat{m}(\mathbf{x}_1)$  is estimated on the attribute values from majority samples.

Stage II: Learning the residual by expanding the c-GAN: Similar to the setting of c-VAE, we now expand the generator  $G_{\theta_{mjr}}$  with several additional layers, denoted as  $G_{\theta_{res}}$ , starting with the feature map produced by  $G_{\theta_{mjr}}$  before the final activation layer. As illustrated in Fig. 2B, we expand the discriminator with additional layers  $D_{\phi_{res}}$  in a mirror of the expanded generator network. Both

 $G_{\theta_{\text{res}}}$  and  $D_{\theta_{\text{res}}}$  are conditioned on the residual  $\gamma = \mathbf{x}_2 - \hat{m}(\mathbf{x}_1)$ . The expanded networks are trained on minority samples using the adversarial loss while freezing Stage-I network weights  $\theta_{\text{mir}}$  and  $\phi_{\text{mir}}$ :

$$\min_{\theta_{\text{res}}} \max_{\phi_{\text{res}}} \left\{ \mathbb{E}_{\mathbf{y}_{\text{mjr}}, \mathbf{x}_{1}, \gamma \sim p_{\text{data}}} \left[ \log D_{\phi_{\text{mjr}}} (D_{\phi_{\text{res}}}(\mathbf{y}_{\text{mnr}}, \gamma), \mathbf{x}_{1}) \right] \right. \\
+ \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z}), \mathbf{x}_{1}, \gamma \sim p_{\text{data}}} \left[ \log \left( 1 - D_{\phi_{\text{mjr}}} (D_{\phi_{\text{res}}}(\hat{\mathbf{y}}_{\text{mnr}}, \gamma), \mathbf{x}_{1}) \right) \right] \right\}$$
where  $\hat{\mathbf{y}}_{\text{mnr}} = G_{\theta_{\text{res}}} (G_{\theta_{\text{mir}} \setminus \sigma}(\mathbf{z}, \mathbf{x}_{1}), \mathbf{x}_{1}, \gamma), \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  (8)

Similarly, with Equation (8), the expanded  $G_{\theta_{\rm res}}$  learns to use the residual attribute  $\gamma$  to change the distribution of generated majority samples to one that aligns with the distribution of real minority images  $\mathbf{y}_{\rm mnr}$ . At the same time, the expanded discriminator  $D_{\phi_{\rm res}}$  is encouraged to change the residual features on the real/generated minority sample in order to leverage the stage-I discriminator  $D_{\phi_{\rm min}}$  that has learned to work with the distribution of majority samples.

One difference between Equation (5) for ProReGen-VAE and Equation (8) for ProReGen-GAN is the distribution over which the sample **z** is taken. This is inherently determined by the training loss of the two models, where the likelihood loss of VAE is calculated over the posterior distribution of **z** conditioned on an observed image (emphasizing instance-level reconstruction) *vs.* the adversarial loss in GAN is calculated over the prior density of **z** (emphasizing distribution-level distance).

#### 5 EXPERIMENTS AND RESULTS

**Data:** We consider Colored-MNIST (Lee et al., 2021), MNIST-Correlation (Mu & Gilmer, 2019), and Corrupted-CIFAR10 (Hendrycks & Dietterich, 2019). We assume known labels of the attributes that are correlated. For the two MNIST datasets, we curated high levels of correlation strengths at 95%, 98%, 99%, and 99.5%, where the % represents the percentage of majority training samples in the training data. For Corrupted-CIFAR10 derived from natural images, we considered less extreme correlation strengths of 70% and 80%. For each dataset, we included a *balanced* version without any attribute correlations to both establish *oracle* attribute classifiers and establish a reference performance for all models considered. Details of the dataset will be described in their respective sections below, and test accuracies of the oracle classifiers for each dataset are presented in Appendix A. Moreover, we present the training data distribution for each dataset in Appendix B.

**Baselines:** We considered c-VAE and c-GAN baselines with the following strategies for mitigating attribute correlations: 1) naive, 2) re-weighting, achieved with upsampling minority samples using the *weighted random sampler* and 3) pseudo-supervision on counterfactual generations, represented by causal-cHVAE (Ribeiro et al., 2023) where a classifier is used to finetune the model in an optional second-stage of counterfactual generation, and causal-GAN (Kocaoglu et al., 2017) where the attribute classifier is incoporated in the end-to-end adversarial training. We present the architectural details of all baselines and ProReGen in Appendix C.

**Evaluation:** We evaluated the performance of all DGMs in generating both majority and minority samples. To generate with the trained DGMs, we sampled from  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and generated a total of 25,000 samples, with equal number of samples for each unique attribute combination, per dataset. We evaluated generated samples using: 1) <u>correctness</u>, measuring the ratio of generations in which the attributes evaluated by the oracle classifier match the intended attributes; 2) <u>Fréchet Inception Distance (FID)</u> (Heusel et al., 2017), measuring the quality and diversity of generations by comparing the representations (retrieved from *InceptionV3* network) of the generated samples against a test set of diverse real samples; and 3) <u>coverage & density</u> (Naeem et al., 2020), measuring the diversity and fidelity, respectively, of generations compared with a test set of diverse real samples.

#### 5.1 EXPERIMENTS & RESULTS ON COLORED-MNIST

Settings of Attribute Correlations: Colored-CMNIST (Lee et al., 2021) is a commonly-used benchmark for synthesizing attribute correlations in the training data. It is an MNIST-variant with a distinct majority color for each of the 10 digits (e.g., orange as a majority color for digit 1). This creates a correlation between digit and color attributes, both discrete but non-binary. All baselines as described in were considered in this dataset, with the exception of causal-GAN which was designed to work with binary labels only in the original paper Kocaoglu et al. (2017). We consider digit  $\rightarrow$  color as the causal direction  $x_1 \rightarrow x_2$  for our experiments.

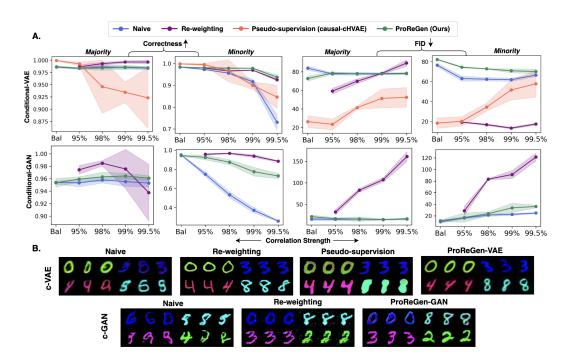


Figure 3: A: Correctness and FID of majority and minority generations from ProReGen vs. baselines for Colored-MNIST. B: Visual examples of minority generations at 99.5% correlation ratio.

**Results and Analysis:** We present some representative quantitative metrics and visual examples on Colored-MNIST in Fig. 3. Complete results are included in Appendix F.1. As shown, ProReGen (green) in general improved the correctness of minority generations in comparison to the naively trained baseline (blue, significantly more in c-GAN), at comparable or slightly worsened quality metrics. In comparison, while causal-cHVAE (red, the baseline that leverages pseudo-supervision) also improved the correctness of minority generations, this improvement was obtained at the expense of degraded correctness in majority generations, suggesting that the use of pseudo-supervision has introduced trade-off in the correctness of majority *vs.* minority generations at higher correlation strengths. ProReGen, in comparison, was relatively stable across correlation strengths for majority generations. Note that the relatively strong quality metrics of causal-cHAVE may be due to its base hierarchical VAE architecture that was different from the rest of the VAE models considered.

Re-weighting resulted in comparable correctness metrics in both majority and minority generations in comparison to ProReGen. However, in both c-VAE and c-GAN, signs of overfitting could be observed in minority generations in the re-weighted baseline (Fig. 3), and hence limited diversity of its generations. In re-weighted c-VAE, the higher coverage of minority generations (Fig. 8 in Appendix F.1) in comparison to ProReGen can seem counterintuitive. A potential reason for it was overfitting to the limited training minority samples and hence being closer to the real data (and sharper in perceptual quality) in re-weighted c-VAE, than the generations from learned distribution in our approach, which were comparatively blur; thus, the coverage value, which was calculated using nearest-neighborhood, resulted in higher value for re-weighted baseline. In re-weighted c-GAN, this led to significant degradation in generation quality, and hence resulting in worse FID, coverage, and density values, with a significant drop, as correlation increase (Fig. 3).

#### 5.2 EXPERIMENTS & RESULTS ON MNIST-CORRELATION

Settings of Attribute Correlations: MNIST-Correlation (Mu & Gilmer, 2019) is another MNIST variant where most of the even digits are clean and most of the odd digits include zigzag, hence resulting in a correlation between attributes  $\mathbf{x_1} = \{\text{even, odd}\}$  and  $\mathbf{x_2} = \{\text{clean, zigzag}\}$ . We similarly created correlation strengths at 95%, 98%, 99%, and 99.5% following (Goel et al., 2020). Along with the information on *presence / absence* of zigzag, we also added the coordinates of the end points of zigzag (mid-point of image in case of *clean* image) as additional feature at-

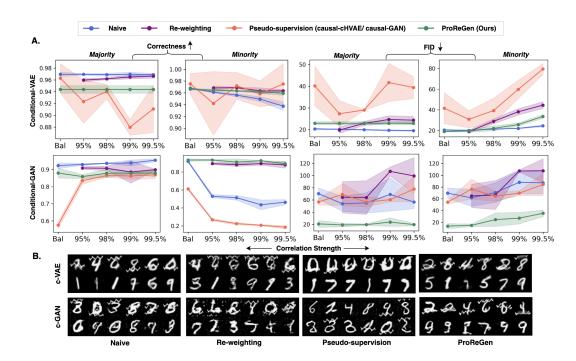


Figure 4: A: Correctness and FID of majority and minority generations from ProReGen vs. baselines for MNIST-Correlation. B: Visual comparison of minority generations at 99.5% correlation ratio. In each image grid, the intended generation is zigzag-even on the top, and clean-odd for the bottom.

tributes in  $\mathbf{x_2}$  to represent residual attributes that cannot be predicted from  $\mathbf{x_1}$  but will contribute to the generation of images. Additional details on this are included in Appendix C. We consider even / odd  $\rightarrow$  presence / absence of zigzag as the causal direction for our experiments.

**Results and Analysis:** We present quantitative results and visual examples for MNIST-Correlation in Fig. 4, with complete results in Appendix F.2. Compared to naive c-VAE (blue), ProReGen-VAE (green) was able to improve the correctness of minority generations, at some degradation in the correctness of majority generations and similar or slightly worsened quality metrics. In comparison, causal-cHVAE (red) was inconsistent in improving the correctness of minority generation at more significance deterioration of both majority correctness and FID metrics. Re-weighting (purple) delivered similar correctness in minority generations and quality metrics in majority generations, but better majority correctness and worsened FID (reflecting diversity issue) in minority generations.

Compared to naively trained c-GAN (blue), ProReGen-GAN (green) significantly improved the correctness of minority generations along with significantly improved FID in both generations with moderate degradation of correctness in majority generations. Causal-GAN (red) was not successful in improving the correctness of minority generations, with FID metrics similar to the naive baseline. Reweighting (purple) improved correctness of minority generations with slight compromise in the correctness of majority generations, but also worsened FID metrics.

#### 5.3 EXPERIMENTS & RESULTS ON CORRUPTED-CIFAR10

Settings of Attribute Correlations: Finally, we adopted CIFAR10 (Krizhevsky et al., 2009) and curated it following the practice in Hendrycks & Dietterich (2019) to create a correlation between object classes and image corruption types. More specifically, we considered five different object classes,  $\mathbf{x_1} = \{\text{car, bird, dog, horse, ship}\}$ , and applied a unique type of corruption,  $\mathbf{x_2} = \{\text{gaussian noise, shot noise, impulse noise, contrast, brightness}\}$ , respectively, to the majority of the training samples per object class. The minority samples have remaining corruptions uniformly sampled at random. We considered less extreme correlation strengths at the level of 70% and 80%.

**Results and Analysis:** We tested only ProReGen-GAN on this dataset due to quality issue (blurring of corruption details) of VAE-based models on these natural images. Representative results in Fig.

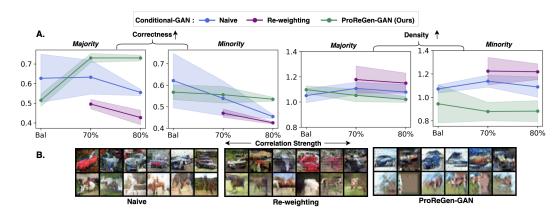


Figure 5: A: Correctness and density of generated images on Corrupted-CIFAR10. B: Visual examples of minority generations at 80% correlation ratio. In each image grid, the intended generation is {car, impulse noise} on the top, and {horse, brightness} for the bottom.

Table 1: A comparison of progressive two-stage training vs. simultaneous training of  $g_{mjr}$  and  $g_{res}$  in ProReGen-GAN for 95% correlation strength in Colored-MNIST

		Correctness	FID	Coverage	Density
Two-Staged Training	Majority	$0.9592 \pm 0.004$	$16.9488 \pm 2.3446$	$0.9003 \pm 0.0275$	$0.7628 \pm 0.0386$
	Minority	$0.9256 \pm 0.0257$	$17.2562 \pm 7.8816$	$0.7519 \pm 0.0903$	$0.6089 \pm 0.1216$
Single-Staged Training	Majority	$0.9289 \pm 0.0369$	$29.0843 \pm 6.3397$	$0.7636 \pm 0.0216$	$0.5585 \pm 0.0114$
	Minority	$0.3557 \pm 0.1689$	$65.0227 \pm 12.3787$	$0.0432 \pm 0.0335$	$0.0320 \pm 0.0220$

5 showed that ProReGen-GAN was able to improve the correctness of both minority and majority generations, although at the expense of degration of image generation qualities.

#### 5.4 Additional Ablation Studies

Effects of Progressive training: To demonstrate the benefit of the progressive two-stage training, we performed an ablation of ProReGEN where the model architecture remained the same but  $g_{\rm mjr}$  and  $g_{\rm res}$  were optimized simultaneous vs. progressively in two stages. As shown in Table 1, without progress training, there was minimal to no effect on the correctness of majority generations but some impact on its quality (e.g., a drop of 27% in density). The generation of minority samples however was significantly worsened (e.g., a drop of 62% in correctness and nearly three times worse in FID). We further demonstrate this with sample minority generations in Fig. 6 in Appendix.

Effect of assumed causal directions between attributes: We examined the effect of the assumption of causal directions between attributes by inverting the causal direction  $digit \rightarrow color$  to  $color \rightarrow digit$  for Colored-MNIST. We considered ProReGen-GAN for our analysis. We observed that the performance with the inverted causal direction  $color \rightarrow digit$  was suboptimal, with only  $0.0811 \pm 0.0127$  correctness, on average, of minority generations vs.  $0.9256 \pm 0.0257$ , on average, with  $digit \rightarrow color$ . The correctness of majority generations was similar. This indicated that learning the residual for digit conversion was much more difficult. We present the generation samples along with additional results in Appendix E. This suggests that the difficulty, and hence performance of residual generation task, is influenced by the causal direction assumed and should be used to design the attribute causal direction for ProReGen in practice (unless the true causal direction is known).

#### 6 CONCLUSION

We present ProReGen, a novel DGM-design that employs progressive training and leverages majority training samples to learn most part of the generation task while employing minority training samples to only learn the residual information. We demonstrate its benefit in improving generation correctness against the baselines using synthetic and natural images at different correlation ratios.

## REPRODUCIBILITY STATEMENT

We present the design of network architectures and training details used in our proposed method in Appendix C. Moreover, we provide reference to the official code repositories employed for experimentation with two of our baselines in the same section. Moreover, we share the training data distribution of the datasets used to present our results in Appendix B.

# REFERENCES

- Jaeju An, Taejune Kim, Donggeun Ko, Sangyup Lee, and Simon S Woo. A^ 2: Adaptive Augmentation for Effectively Mitigating Dataset Bias. In Proceedings of the Asian Conference on Computer Vision, pp. 4077–4092, 2022.
- Joey Bose, Ricardo Pio Monti, and Aditya Grover. Controllable generative modeling via causal reasoning. Transactions on Machine Learning Research, 2022.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters, 2018.
- Karan Goel, Albert Gu, Yixuan Li, and Christopher Re. Model patching: Closing the subgroup performance gap with data augmentation. In <u>International Conference on Learning Representations</u>, 2020.
- Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen. Attgan: Facial attribute editing by only changing what you want. <u>IEEE transactions on image processing</u>, 28(11):5464–5478, 2019.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. arXiv preprint arXiv:1903.12261, 2019.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. <u>Advances in</u> neural information processing systems, 30, 2017.
- Guillaume Jeanneret, Loïc Simon, and Frédéric Jurie. Diffusion models for counterfactual explanations. In Proceedings of the Asian Conference on Computer Vision, pp. 858–876, 2022.
- Eungyeup Kim, Jihyeon Lee, and Jaegul Choo. Biaswap: Removing dataset bias with bias-tailored swapping augmentation. In <u>Proceedings of the IEEE/CVF International Conference on Computer Vision</u>, pp. 14992–15001, 2021.
- Murat Kocaoglu, Christopher Snyder, Alexandros G Dimakis, and Sriram Vishwanath. Causalgan: Learning causal implicit generative models with adversarial training. <a href="mailto:arXiv:1709.02023">arXiv:1709.02023</a>, 2017.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Jungsoo Lee, Eungyeup Kim, Juyoung Lee, Jihyeon Lee, and Jaegul Choo. Learning debiased representation via disentangled feature augmentation. <u>Advances in Neural Information Processing Systems</u>, 34:25123–25133, 2021.
- Miguel Monteiro, Fabio De Sousa Ribeiro, Nick Pawlowski, Daniel C Castro, and Ben Glocker. Measuring axiomatic soundness of counterfactual image models. In <a href="The Eleventh International Conference">The Eleventh International Conference on Learning Representations, 2022.</a>
- Norman Mu and Justin Gilmer. Mnist-c: A robustness benchmark for computer vision. <u>arXiv</u> preprint arXiv:1906.02337, 2019.
  - Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. In <u>International conference on machine learning</u>, pp. 7176–7185. PMLR, 2020.

- Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects. Biometrika, 108(2):299–319, 2021.
- Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei Efros, and Richard Zhang. Swapping autoencoder for deep image manipulation. <u>Advances in Neural Information</u> Processing Systems, 33:7198–7211, 2020.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv preprint arXiv:1511.06434, 2015.
- Fabio De Sousa Ribeiro, Tian Xia, Miguel Monteiro, Nick Pawlowski, and Ben Glocker. High fidelity image counterfactuals with probabilistic causal models. In <u>International Conference on Machine Learning</u>, pp. 7390–7425. PMLR, 2023.
- Peter M Robinson. Root-n-consistent semiparametric regression. <u>Econometrica</u>: journal of the Econometric Society, pp. 931–954, 1988.
- Pau Rodríguez, Massimo Caccia, Alexandre Lacoste, Lee Zamparo, Issam Laradji, Laurent Charlin, and David Vazquez. Beyond Trivial Counterfactual Explanations With Diverse Valuable Explanations. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 1056–1065, October 2021.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally Robust Neural Networks. In International Conference on Learning Representations, 2019.
- Axel Sauer and Andreas Geiger. Counterfactual Generative Networks. In <u>International Conference</u> on Learning Representations, 2020.
- Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. Advances in neural information processing systems, 28, 2015.
- Frederik Träuble, Elliot Creager, Niki Kilbertus, Francesco Locatello, Andrea Dittadi, Anirudh Goyal, Bernhard Schölkopf, and Stefan Bauer. On disentangled representations learned from correlated data. In International conference on machine learning, pp. 10401–10412. PMLR, 2021.
- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2097–2106, 2017.
- Wenqian Ye, Guangtao Zheng, Xu Cao, Yunsheng Ma, and Aidong Zhang. Spurious correlations in machine learning: A survey. arXiv preprint arXiv:2402.12715, 2024.

Table 2: Test Accuracy of Oracle Classifiers for Colored-MNIST, MNIST-Correlation, and Corrupted-CIFAR10

	Oracle Classifier For	Average Test Accuracy
Colored-MNIST	Digit	0.958
Color cu-Minist	Color	1.0
MNIST-Correlation	Even/Odd Digit Type	0.965
WINDI-Correlation	Presence/Absence of Zigzag	1.0
Corrupted-CIFAR10	Object Type	0.849
Corrupted-Cirricio	Corruption Type	0.99

Table 3: Counts of Majority and Minority Samples Across Varying Levels of Correlation Strengths Explored for Colored-MNIST. For the *Balanced* setting, we consider equal number of samples per (digit, color) combination, *i.e.*, around 550 samples per combination. The total number of unique combinations is 100.

<b>Correlation Strength</b>	Minority	Majority
95%	2450	52552
98%	986	54014
99%	492	54510
99 5%	249	54751

Table 4: Counts of Majority and Minority Samples Across Varying Levels of Correlation Strengths Explored for MNIST-Correlation. For the *Balanced* setting, we consider equal number of samples per attribute combination, *i.e.*, around 20,000 samples per combination. The total number of unique combinations is 4.

Correlation Strength	Minority	Majority
95%	2104	40000
98%	816	40000
99%	404	40000
99.5%	200	40000

Table 5: Counts of Majority and Minority Samples Across Varying Levels of Correlation Strengths Explored for Corrupted-CIFAR10. For the *Balanced* setting, we consider equal number of samples per attribute combination, *i.e.*, around 800 samples per combination. The total number of unique combinations is 25.

Correlation Strength	Minority	Majority
70%	6000	14000
80%	4000	16000

# A ORACLE PERFORMANCE FOR EACH DATASET

We present the performance oracle classifiers for Colored-MNIST, MNIST-Correlation, and Corrupted-CIFAR10 in Table 2.

## B TRAINING DATA DISTRIBUTION FOR EACH DATASET

We present the count of majority and minority training samples employed across each correlation ratio for Colored-MNIST in Table 3, for MNIST-Correlation in Table 4, and for Corrupted-CIFAR10 in Table 5.

Table 6: The architecture of expanded generator network for Colored-MNIST. Encoder or discriminator is a mirror of it. Here  $C_{img}$  denotes the number of image channels and  $d_t$  is the dimension of

Part	Output Shape	Layer Information
Input	$(B, C_{img} + d_t, H, W)$	-
Conv Block 2	(B, 64, H, W) (B, 32, H, W) (B, C <sub>img</sub> , H, W)	Conv2d( $C_{\text{img}} + d_t$ , 64, 3, 1, 1), GroupNorm(8, 64), ReLU Conv2d(64, 32, 3, 1, 1), GroupNorm(8, 32), ReLU Conv2d(32, $C_{\text{img}}$ , 3, 1, 1)

Table 7: The architecture of expanded generator network for MNIST-Correlation. Encoder or discriminator is a mirror of it. Here  $C_{img}$  denotes the number of image channels and  $d_t$  is the dimension of  $\gamma$ .

Part	Output Shape	Layer Information
Input	$(B, C_{img} + d_t, H, W)$	-
Conv Block 1 Conv Block 2 Conv Block 3 Conv Block 4 Output Layer	(B, 64, H, W) (B, 64, H, W) (B, 32, H, W) (B, 32, H, W) (B, C <sub>img</sub> , H, W)	$ \begin{array}{l} {\rm Conv2d}(C_{\rm img}+d_t,64,3,1,1),{\rm GroupNorm}(8,64),{\rm ReLU} \\ {\rm Conv2d}(64,64,3,1,1),{\rm GroupNorm}(8,64) \\ {\rm Conv2d}(64,32,3,1,1),{\rm GroupNorm}(8,32),{\rm ReLU} \\ {\rm Conv2d}(32,32,3,1,1),{\rm GroupNorm}(8,32) \\ {\rm Conv2d}(32,C_{\rm img},3,1,1) \end{array} $

# C IMPLEMENTATION DETAILS

**Expanded Network Architecture.** We employ a stack of convolution and transposed convolution layers to implement the c-VAE (Sohn et al., 2015) and follow DCGAN-like architectural setup for c-GAN (Radford et al., 2015). For causal-cHVAE and causal-GAN, we follow their official code repositories: Ribeiro et al. (2023) for causal-cHVAE and Kocaoglu et al. (2017) for causal-GAN.

In both c-VAE and c-GAN, since the residual effect generator only requires adjusting residual features on images for which most of the generative factors have already been produced by the majority DGM, the expanded layers are lightweight, comprising relatively low parameter count than the majority DGMs. Moreover, they are also designed such that they maintain the spatial dimension of the input image (e.g., using convolutional operations with kernel = 3, stride = 1, padding = 1). We present the details of their network architecture in Table 6-8.

Conditional Information  $\gamma$ . Residual (orthogonal) attribute  $\gamma=x_2-m(x_1)$  is broadcasted to match spatial size of the input and concatenated as additional channels to provide it as conditional information to the expanded layers. When we have high-dimensional attribute  $x_2$ , where only part of it is predictable from  $x_1$ , we predict the predictable dimensions by  $m(x_1)$  to get  $\gamma$ . As implementation choice, we can then either keep the additional (unpredictable) dimensions as additional channels or use them to manipulate (e.g., mask) predictable dimensions. We employ the latter approach for MNIST-Correlation for ProReGen, where we mask the  $\gamma$  information using the information of line joining the two end-points of zigzag. We employ the same approach in the encoder of naively-trained conditional-VAE and its re-weighted version. However, employing such masking in naively-trained conditional-GAN led to the model ignoring the even/odd label information, potentially due to the discriminator relying on the now easier feature, zigzag (due to its masking). Therefore, we appended the coordinate information as additional channels during conditioning in conditional-GAN.

Table 8: The architecture of expanded generator network for Corrupted-CIFAR10. The discriminator network is a mirror of it. Here  $C_{img}$  denotes the number of image channels and  $d_t$  denotes the dimension of  $\gamma$  and  $\mathbf{x}_1$ , which are the same value.

Part	Output Shape	Layer Information
Residual Block 1	$(B, C_{mid}, H, W)$	Concatenate $h_{\mathrm{mjr}}$ , $\gamma$ , $\mathbf{x}_1$ along channel dim $\mathrm{Conv2d}(C_{\mathrm{img}} + 2d_t, C_{\mathrm{mid}}, 3, 1, 1)$ , $\mathrm{GroupNorm}(8, C_{\mathrm{mid}})$ , $\mathrm{ReLU}$ $\mathrm{Conv2d}(C_{\mathrm{mid}}, C_{\mathrm{mid}}, 3, 1, 1)$ , $\mathrm{GroupNorm}(8, C_{\mathrm{mid}})$ $\mathrm{Skip}$ connection: $\mathrm{Conv2d}(C_{\mathrm{img}}, C_{\mathrm{mid}}, 1, 1, 0)$ (or Identity if channels match) $\mathrm{Element\text{-}wise}$ addition (residual connection)
Residual Block 2	$(B, C_{mid2}, H, W)$	Concatenate Residual Block 1 output, $\gamma$ , $\mathbf{x}_1$ along channel dim Conv2d( $C_{\text{mid}}+2d_t, C_{\text{mid2}}, 3, 1, 1$ ), GroupNorm(8, $C_{\text{mid2}}$ ), ReLU Conv2d( $C_{\text{mid2}}, C_{\text{mid2}}, 3, 1, 1$ ), GroupNorm(8, $C_{\text{mid2}}$ ) Skip connection: Conv2d( $C_{\text{mid}}, C_{\text{mid2}}, 1, 1, 0$ ) (or Identity if channels match) Element-wise addition (residual connection)
Output Layer	$(B, C_{img}, H, W)$	$Conv2d(C_{mid2}, C_{img}, 3, 1, 1)$

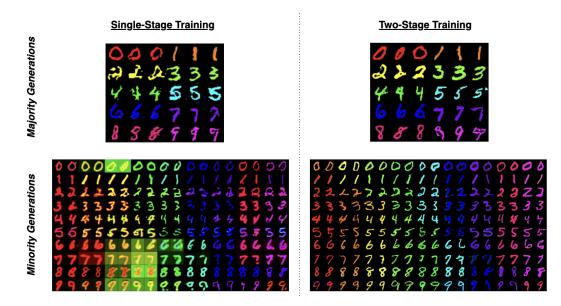


Figure 6: Sample majority and minority generations with two-stage training of  $g_{mjr}$  and  $g_{res}$  vs. single-stage training for 95% correlation ratio in Colored-MNIST.

# D EFFECT OF TWO-STAGE TRAINING

Two-stage training of  $g_{\rm mjr}$  and  $g_{\rm res}$  instead of single-stage training is beneficial for the quality of majority generations and overall success (correctness, FID, coverage, and density) of minority generations. We present visual examples of generations in Fig. 6.

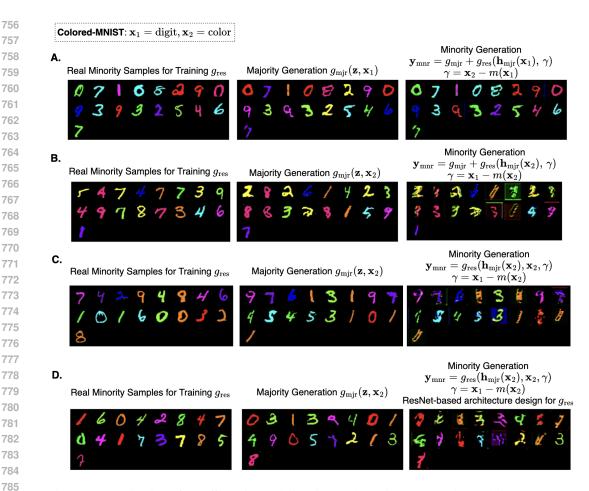


Figure 7: Examination of the effect of causal direction on the residual generation task for ProReGen-GAN, considering 95% correlation ratio in Colored-MNIST. Real minority samples used for training  $g_{res}$  are shown in *left* with corresponding majority (middle) and minority (right) generation samples. We employ the causal direction  $digit \rightarrow color$  for Colored-MNIST in our main experiments and present the sample results in **A**. We explore the effect of inverting the causal direction to  $color \rightarrow digit$  in **B**, **C**, and **D**, where the roles of  $\mathbf{x_1}$  and  $\mathbf{x_2}$  are reversed.

#### E EFFECT OF ASSUMPTION OF CAUSAL DIRECTION

We present the comparison of generated majority and corresponding minority samples when considering causal direction  $x_1 \to x_2$  vs.  $x_2 \to x_1$  in Fig. 7. We consider ProReGen-GAN to present our analysis.

In Fig. 7-B, we simply invert the causal direction, while employing the same additive formulation for  $y_{mnr}$  as in our main experiments and keeping the network architecture style for  $g_{res}$  consistent.

We further experiment with the general formulation  $\mathbf{y}_{mnr} = g_{res}(\mathbf{h}_{mjr}(\mathbf{x}_1), \mathbf{x}_1, \gamma)$  using: 1) the same network architecture design for  $g_{res}$  as in Fig. 7-B, and 2) ResNet-style architecture design for  $g_{res}$  such that the residual operation occurs implicitly within the network, to assess their potential benefit for the residual generation task with the inverted causal direction  $color \rightarrow digit$ . However, no noticeable improvement was observed as shown in Fig. 7C-D.

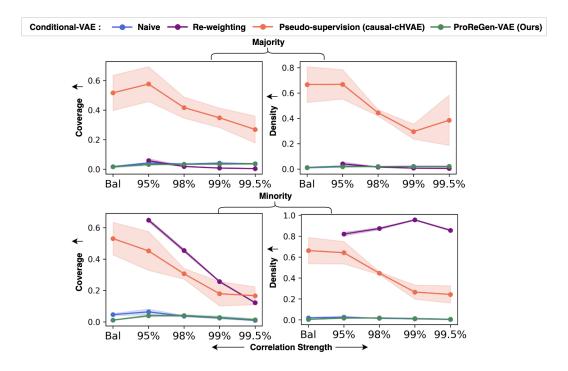


Figure 8: Comparison of coverage and density metric values of ProReGen-VAE against the baselines for Colored-MNIST dataset.

#### F ADDITIONAL RESULTS

# F.1 COLORED-MNIST

We present the comparison of coverage and density metric values of ProReGen against the baselines in Fig. 8 and Fig. 9.

#### F.2 MNIST-CORRELATION

We present the comparison of coverage and density metric values of ProReGen against the baselines for MNIST-Correlation in Fig. 10 and Fig. 11.

## F.3 CORRUPTED-CIFAR10

We present the comparison of coverage metric value of ProReGen against the baselines for Corrupted-CIFAR10 in Fig. 12.

#### G LLM USAGE

We used the LLM tool, ChatGPT, at limited capacity. ChatGPT was leveraged for improving the quality of sentences to provide better readability and for grammatical corrections. Moreover, we utilized it to generate some portions of the graph creation scripts.

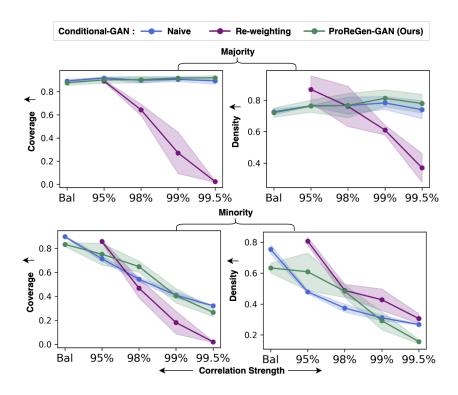


Figure 9: Comparison of coverage and density metric values of ProReGen-GAN against the baselines for Colored-MNIST dataset.

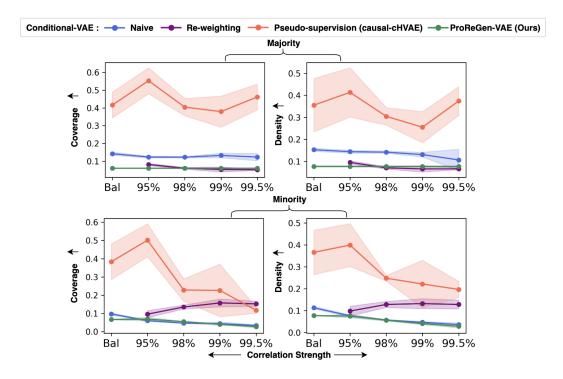


Figure 10: Comparison of coverage and density metric values of ProReGen-VAE against the baselines for MNIST-Correlation dataset.

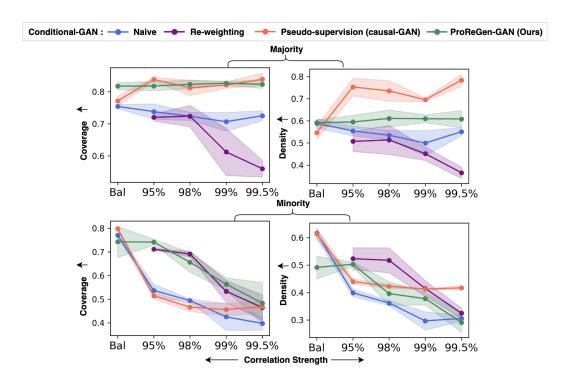


Figure 11: Comparison of coverage and density metric values of ProReGen-GAN against the baselines for MNIST-Correlation dataset.

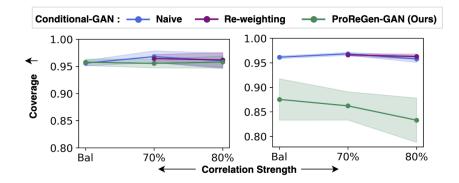


Figure 12: Comparison of coverage metric value of ProReGen-GAN against the baselines for Corrupted-CIFAR10 dataset.