

---

# D2C: Diffusion-Decoding Models for Few-Shot Conditional Generation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Conditional generative models of high-dimensional images have many applications,  
2 but supervision signals from conditions to images can be expensive to acquire.  
3 This paper describes Diffusion-Decoding models with Contrastive representations  
4 (D2C), a paradigm for training unconditional variational autoencoders (VAEs)  
5 for few-shot conditional image generation. D2C uses a learned diffusion-based  
6 prior over the latent representations to improve generation and contrastive self-  
7 supervised learning to improve representation quality. D2C can adapt to novel  
8 generation tasks conditioned on labels or manipulation constraints, by learning from  
9 as few as 100 labeled examples. On conditional generation from new labels, D2C  
10 achieves superior performance over state-of-the-art VAEs and diffusion models.  
11 On conditional image manipulation, D2C generations are two orders of magnitude  
12 faster to produce over StyleGAN2 ones and are preferred by 50% – 60% of the  
13 human evaluators in a double-blind study.

## 14 1 Introduction

15 Generative models trained on large amounts of unlabeled data have achieved great success in various  
16 domains including images [8, 47, 72, 40], text [53, 2], audio [24, 68, 88, 59], and graphs [34, 64].  
17 However, downstream applications of generative models are often based on various conditioning  
18 signals, such as labels [58], text descriptions [57], reward values [96], or similarity with existing  
19 data [43]. While it is possible to directly train conditional models, this often requires large amounts  
20 of paired data [54, 71] that are costly to acquire. Hence, it would be desirable to learn conditional  
21 generative models using large amounts of unlabeled data and as little paired data as possible.

22 Contrastive self-supervised learning (SSL) methods can greatly reduce the need for labeled data in  
23 discriminative tasks by learning effective representations from unlabeled data [90, 35, 33], and have  
24 also been shown to improve few-shot learning [37]. It is therefore natural to ask if they can also  
25 be used to improve few-shot generation. Latent variable generative models (LVGM) are a natural  
26 candidate for this, since they already involve a low-dimensional, structured latent representation  
27 of the data they generate. However, popular LVGMs, such as generative adversarial networks  
28 (GANs, [32, 47]) and diffusion models [40, 80], lack explicit tractable functions to map inputs to  
29 representations, making it difficult to optimize latent variables with SSL. Variational autoencoders  
30 (VAEs, [49, 74]), on the other hand, can naturally adopt SSL through their encoder model, but they  
31 typically have worse sample quality.

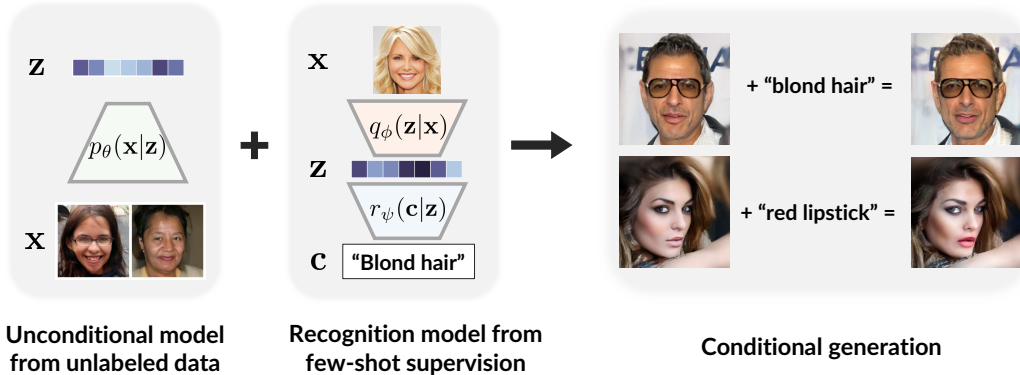


Figure 1: Few-shot conditional generation with the unconditional D2C model (left). With a recognition model over the latent space (middle), D2C can generate samples for novel conditions, such as image manipulation (right). These conditions can be defined with very few labels.

32 In this paper, we propose Diffusion-Decoding models with Contrastive representations (D2C), a  
 33 special VAE that is suitable for conditional few-shot generation. D2C uses contrastive self-supervised  
 34 learning methods to obtain a latent space that inherits the transferrability and few-shot capabilities of  
 35 self-supervised representations. Unlike other VAEs, D2C learns a diffusion model over the latent  
 36 representations. This latent diffusion model ensures that D2C uses the same latent distribution for  
 37 both training and generation. We provide a formal argument to explain why this approach may lead  
 38 to better sample quality than existing hierarchical VAEs. We further discuss how to apply D2C to  
 39 few-shot conditional generation where the conditions are defined through labeled examples and/or  
 40 manipulation constraints. Our approach combines a discriminative model providing conditioning  
 41 signal and generative diffusion model over the latent space, and is computationally more efficient  
 42 than methods that act directly over the image space (Figure 1).

43 We evaluate and compare D2C with several state-of-the-art generative models over 6 datasets. On  
 44 unconditional generation, D2C outperforms state-of-the-art VAEs and is competitive with diffusion  
 45 models under similar computational budgets. On conditional generation with 100 labeled examples,  
 46 D2C significantly outperforms state-of-the-art VAE [87] and diffusion models [80]. D2C can also  
 47 learn to perform certain image manipulation tasks from as few as 100 labeled examples. Notably, for  
 48 manipulating images, D2C is two orders of magnitude faster than StyleGAN2 [101] and preferred by  
 49 50% – 60% of human evaluations, which to our best knowledge is the first for any VAE model.

## 50 2 Background

51 **Latent variable generative models** A latent variable generative model (LVGM) is posed as a con-  
 52 ditional distribution  $p_\theta : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{X})$  from a latent variable  $\mathbf{z}$  to a generated sample  $\mathbf{x}$ , parametrized  
 53 by  $\theta$ . To acquire new samples, LVGMs draw random latent variables  $\mathbf{z}$  from some distribution  
 54  $p(\mathbf{z})$  and map them to image samples through  $\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{z})$ . Most LVGMs are built on top of three  
 55 paradigms: variational autoencoders (VAEs, [49, 74]), Normalizing Flows (NFs, [26, 27]), Generative  
 56 Adversarial Networks (GANs, [32]), and diffusion / score-based generative models [40, 81].

57 Particularly, VAEs use an inference model from  $\mathbf{x}$  to  $\mathbf{z}$  for training. Denoting the inference distribution  
 58 from  $\mathbf{x}$  to  $\mathbf{z}$  as  $q_\phi(\mathbf{z}|\mathbf{x})$ , the generative distribution from  $\mathbf{z}$  to  $\mathbf{x}$  as  $p_\theta(\mathbf{x}|\mathbf{z})$ , VAEs are trained by  
 59 minimizing the following upper bound of negative log-likelihood:

$$L_{\text{VAE}} = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [-\log p(\mathbf{x}|\mathbf{z})] + D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))] \quad (1)$$

60 where  $p_{\text{data}}$  is the data distribution and  $D_{\text{KL}}$  is the KL-divergence.

61 **Diffusion models** Diffusion models [78, 40, 80] produce samples by reversing a Gaussian diffusion  
 62 process. We use the index  $\alpha \in [0, 1]$  to denote the particular noise level of a noisy observation  
 63  $\mathbf{x}^{(\alpha)} = \sqrt{\alpha}\mathbf{x} + \sqrt{1-\alpha}\epsilon$ , where  $\mathbf{x}$  is the clean observation and  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$  is a standard Gaussian  
 64 distribution; as  $\alpha \rightarrow 0$ , the distribution of  $\mathbf{x}^{(\alpha)}$  converges to  $\mathcal{N}(0, \mathbf{I})$ . Diffusion models are typically

65 parametrized as reverse noise models  $\epsilon_\theta(\mathbf{x}^{(\alpha)}, \alpha)$  that predict the noise component of  $\mathbf{x}^{(\alpha)}$  given  
 66 a noise level  $\alpha$ , and trained to minimize  $\|\epsilon_\theta(\mathbf{x}^{(\alpha)}, \alpha) - \epsilon\|_2^2$ , the mean squared error loss between  
 67 the true noise and predicted noise. Given any non-decreasing series  $\{\alpha_i\}_{i=0}^T$  between 0 and 1, the  
 68 diffusion objective for a clean sample from the data  $\mathbf{x}$  is:

$$\ell_{\text{diff}}(\mathbf{x}; w, \theta) := \sum_{i=0}^T w(\alpha_i) \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\epsilon - \epsilon_\theta(\mathbf{x}^{(\alpha_i)}, \alpha_i)\|_2^2], \quad \mathbf{x}^{(\alpha_i)} := \sqrt{\alpha_i} \mathbf{x} + \sqrt{1 - \alpha_i} \epsilon \quad (2)$$

69 where  $w : \{\alpha_i\}_{i=0}^T \rightarrow \mathbb{R}_+$  controls the loss weights for each  $\alpha$ . When  $w(\alpha) = 1$  for all  $\alpha$ , we recover  
 70 the denoising score matching objective for training score-based generative models [81].

71 Given an initial sample  $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I})$ , diffusion models acquires clean samples (*i.e.*, samples of  $\mathbf{x}_1$ )  
 72 through a gradual denoising process, where samples with reducing noise levels  $\alpha$  are produced (*e.g.*,  
 73  $\mathbf{x}_0 \rightarrow \mathbf{x}_{0.3} \rightarrow \mathbf{x}_{0.7} \rightarrow \mathbf{x}_1$ ). In particular, Denoising Diffusion Implicit Models (DDIMs, [80]) uses  
 74 an Euler discretization of some neural ODE [13] to produce samples (Figure 2, left).

75 We provide a more detailed description for training diffusion models in Appendix A.1 and sampling  
 76 from DDIM in Appendix A.2. For conciseness, we use the notation  $p^{(\alpha)}(\mathbf{x}^{(\alpha)})$  to denote the marginal  
 77 distribution of  $\mathbf{x}^{(\alpha)}$  under the diffusion model, and  $p^{(\alpha_1, \alpha_2)}(\mathbf{x}^{(\alpha_2)} | \mathbf{x}^{(\alpha_1)})$  to denote the diffusion  
 78 sampling process from  $\mathbf{x}^{(\alpha_1)}$  to  $\mathbf{x}^{(\alpha_2)}$  (assuming  $\alpha_1 < \alpha_2$ ). This notation abstracts away the exact  
 79 sampling procedure of the diffusion model, which depends on choices of  $\alpha$ .

80 **Self-supervised learning of representations** In self-supervised learning (SSL), representations  
 81 are learned by completing certain pretext tasks that do not require extra manual labeling [65, 23]; these  
 82 representations can then be applied to other downstream tasks, often in few-shot or zero-shot scenarios.  
 83 In particular, contrastive representation learning encourages representations to be closer between  
 84 “positive” pairs and further between “negative” pairs; contrastive predictive coding (CPC, [90]), based  
 85 on multi-class classification, have been commonly used in state-of-the-art methods [35, 15, 17, 14, 79].  
 86 Other non-contrastive methods exist, such as BYOL [33] and SimSiam [16], but they usually require  
 87 additional care to prevent the representation network from collapsing.

### 88 3 Problem Statement

89 **Few-shot conditional generation** Our goal is to learn an unconditional generative model  $p_\theta(\mathbf{x})$   
 90 such that it is suitable for conditional generation. Let  $\mathcal{C}(\mathbf{x}, \mathbf{c}, f)$  describe an event that “ $f(\mathbf{x}) = \mathbf{c}$ ”,  
 91 where  $\mathbf{c}$  is a property value and  $f(\mathbf{x})$  is a property function that is unknown at training. In conditional  
 92 generation, our goal is to sample  $\mathbf{x}$  such that the event  $\mathcal{C}(\mathbf{x}, \mathbf{c}, f)$  occurs for a chosen  $\mathbf{c}$ . If we have  
 93 access to some “ground-truth” model that gives us  $p(\mathcal{C}|\mathbf{x}) := p(f(\mathbf{x}) = \mathbf{c}|\mathbf{x})$ , then the conditional  
 94 model can be derived from Bayes’ rule:  $p_\theta(\mathbf{x}|\mathcal{C}) \propto p(\mathcal{C}|\mathbf{x})p_\theta(\mathbf{x})$ . These properties  $\mathbf{c}$  include (but are  
 95 not limited to<sup>1</sup>) labels [58], text descriptions [57, 73], noisy or partial observations [11, 5, 44, 22],  
 96 and manipulation constraints [66]. In many cases, we do not have direct access to the true  $f(\mathbf{x})$ , so  
 97 we need to learn an accurate model from labeled data [6] (*e.g.*,  $(\mathbf{c}, \mathbf{x})$  pairs).

98 **Desiderata** Many existing methods are optimized for some known condition (*e.g.*, labels in con-  
 99 ditional GANs [8]) or assume abundant pairs between images and conditions that can be used for  
 100 pretraining (*e.g.*, DALL-E [71] and CLIP [70] over image-text pairs). Neither is the case in this paper,  
 101 as we do not expect to train over paired data.

102 While high-quality latent representations are not essential to unconditional image generation (*e.g.*,  
 103 autoregressive [89], energy-based [29], and some diffusion models [40]), they can be beneficial  
 104 when we wish to specify certain conditions with limited supervision signals, similar to how SSL  
 105 representations can reduce labeling efforts in downstream tasks. A compelling use case is detecting  
 106 and removing biases in datasets via image manipulation, where we should not only address well-  
 107 known biases a-priori but also address other hard-to-anticipate biases, adapting to societal needs [62].

108 Therefore, a desirable generative model should not only have high sample quality but also contain  
 109 informative latent representations. While VAEs are ideal for learning rich latent representations due  
 110 to being able to incorporate SSL within the encoder, they generally do not achieve the same level of  
 111 sample quality as GANs and diffusion models.

<sup>1</sup>When  $\mathcal{C}$  refers to an event that is always true, we recover unconditioned generation.

Table 1: A comparison of several common paradigms for generative modeling. [Explicit  $\mathbf{x} \rightarrow \mathbf{z}$ ]: the mapping from  $\mathbf{x}$  to  $\mathbf{z}$  is directly trainable, which enables SSL; [No prior hole]: latent distributions used for generation and training are identical (Sec. 4.2), which improves generation; [Non-adversarial]: training procedure does not involve adversarial optimization, which improves training stability.

Model family	Explicit $\mathbf{x} \rightarrow \mathbf{z}$ (Enables SSL)	No prior hole (Better generation)	Non-Adversarial (Stable training)
VAE [49, 74], NF [26]	✓	✗	✓
GAN [32]	✗	✓	✗
BiGAN [28, 30]	✓	✓	✗
DDIM [80]	✗	✓	✓
<b>D2C</b>	✓	✓	✓

## 112 4 Diffusion-Decoding Generative Models with Contrastive Learning

113 To address the above issue, we present Diffusion-Decoding generative models with Contrastive Learn-  
 114 ing (D2C), an extension to VAEs with high-quality samples and high-quality latent representations,  
 115 and are thus well suited to few-shot conditional generation. Moreover, unlike GAN-based methods,  
 116 D2C does not involve unstable adversarial training (Table 1).

117 As its name suggests, the generative model for D2C has two components – *diffusion* and *decoding*;  
 118 the *diffusion* component operates over the latent space and the *decoding* component maps from  
 119 latent representations to images. Let us use the  $\alpha$  index notation for diffusion random variables:  
 120  $\mathbf{z}^{(0)} \sim p^{(0)}(\mathbf{z}^{(0)}) := \mathcal{N}(0, \mathbf{I})$  is the “noisy” latent variable with  $\alpha = 0$ , and  $\mathbf{z}^{(1)}$  is the “clean” latent  
 121 variable with  $\alpha = 1$ . The generative process of D2C, which we denote  $p_\theta(\mathbf{x}|\mathbf{z}^{(0)})$ , is then defined as:

$$122 \mathbf{z}^{(0)} \sim p^{(0)}(\mathbf{z}^{(0)}), \quad \mathbf{z}^{(1)} \sim \underbrace{p_\theta^{(0,1)}(\mathbf{z}^{(1)}|\mathbf{z}^{(0)})}_{\text{diffusion}}, \quad \mathbf{x} \sim \underbrace{p_\theta(\mathbf{x}|\mathbf{z}^{(1)})}_{\text{decoding}}, \quad (3)$$

122 where  $p^{(0)}(\mathbf{z}^{(0)}) = \mathcal{N}(0, \mathbf{I})$  is the prior distribution for the diffusion model,  $p_\theta^{(0,1)}(\mathbf{z}^{(1)}|\mathbf{z}^{(0)})$  is the  
 123 diffusion process from  $\mathbf{z}^{(0)}$  to  $\mathbf{z}^{(1)}$ , and  $p_\theta(\mathbf{x}|\mathbf{z}^{(1)})$  is the decoder from  $\mathbf{z}^{(1)}$  to  $\mathbf{x}$ . Intuitively, D2C  
 124 models produce samples by drawing  $\mathbf{z}^{(1)}$  from a diffusion process and then decoding  $\mathbf{x}$  from  $\mathbf{z}^{(1)}$ .

125 In order to train a D2C model, we use an inference model  $q_\phi(\mathbf{z}^{(1)}|\mathbf{x})$  that predicts proper  $\mathbf{z}^{(1)}$  latent  
 126 variables from  $\mathbf{x}$ ; this can directly incorporate SSL methods [94], leading to the following objective:

$$127 L_{\text{D2C}}(\theta, \phi; w) := L_{\text{D2}}(\theta, \phi; w) + \lambda L_{\text{C}}(q_\phi), \quad (4)$$

$$128 L_{\text{D2}}(\theta, \phi; w) := \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}, \mathbf{z}^{(1)} \sim q_\phi(\mathbf{z}^{(1)}|\mathbf{x})} [-\log p(\mathbf{x}|\mathbf{z}^{(1)}) + \ell_{\text{diff}}(\mathbf{z}^{(1)}; w, \theta)], \quad (5)$$

127 where  $\ell_{\text{diff}}$  is defined as in Eq.(2),  $L_{\text{C}}(q_\phi)$  denotes any contrastive predictive coding objective [90]  
 128 with rich data augmentations [35, 15, 17, 14, 79] (details in Appendix A.3) and  $\lambda > 0$  is a  
 129 weight hyperparameter. The first two terms, which we call  $L_{\text{D2}}$ , contains a “reconstruction loss”  
 130 ( $-\log p(\mathbf{x}|\mathbf{z}^{(1)})$ ) and a “diffusion loss” over samples of  $\mathbf{z}^{(1)} \sim q_\phi(\mathbf{z}^{(1)}|\mathbf{x})$ . We illustrate the D2C  
 131 generative and inference models in Figure 2, and its training procedure in Appendix A.4.

### 132 4.1 Relationship to maximum likelihood

133 The D2 objective ( $L_{\text{D2}}$ ) appears similar to the original VAE objective ( $L_{\text{VAE}}$ ). Here, we make an  
 134 informal statement that the D2 objective function is deeply connected to the variational lower bound  
 135 of log-likelihood; we present the full statement and proof in Appendix B.1.

136 **Theorem 1** (informal). *For any valid  $\{\alpha_i\}_{i=0}^T$ , there exists some weights  $\hat{w} : \{\alpha_i\}_{i=0}^T \rightarrow \mathbb{R}_+$  for the  
 137 diffusion objective such that  $-L_{\text{D2}}$  is a variational lower bound to the log-likelihood, i.e.,*

$$138 -L_{\text{D2}}(\theta, \phi; \hat{w}) \leq \mathbb{E}_{p_{\text{data}}} [\log p_\theta(\mathbf{x})], \quad (6)$$

138 where  $p_\theta(\mathbf{x}) := \mathbb{E}_{\mathbf{z}^{(0)} \sim p^{(0)}(\mathbf{z}^{(0)})} [p_\theta(\mathbf{x}|\mathbf{z}^{(0)})]$  is the marginal probability of  $\mathbf{x}$  under the D2C model.

139 *Proof.* (sketch) The diffusion term  $\ell_{\text{diff}}$  upper bounds the KL divergence between  $q_\phi(\mathbf{z}_1|\mathbf{x})$  and  
 140  $p_\theta^{(1)}(\mathbf{z}^{(1)})$  for suitable weights [40, 80], which recovers a VAE objective.  $\square$

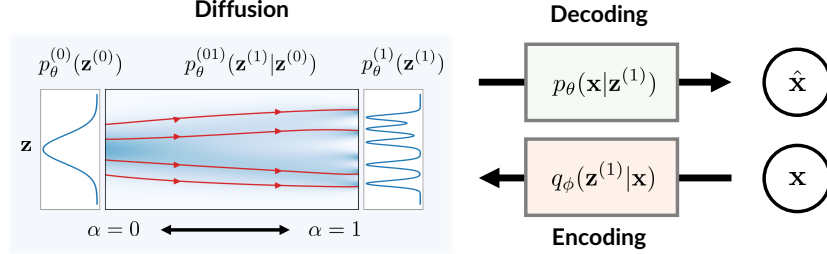


Figure 2: Illustration of components of a D2 model. On top of the encoding and decoding between  $\mathbf{x}$  and  $\mathbf{z}^{(1)}$ , we use a diffusion model to generate  $\mathbf{z}^{(1)}$  from a Gaussian  $\mathbf{z}^{(0)}$ . The red lines describe several smooth ODE trajectories from  $\alpha = 0$  to  $\alpha = 1$  corresponding to DDIM.

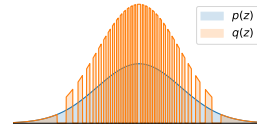
141 **4.2 D2 models address latent posterior mismatch in VAEs**

142 While D2C is a special case of VAE, we argue that D2C is non-trivial in the sense that it addresses a  
 143 long-standing problem in VAE methods [86, 83], namely the mismatch between the prior distribution  
 144  $p_\theta(\mathbf{z})$  and the aggregate (approximate) posterior distribution  $q_\phi(\mathbf{z}) := \mathbb{E}_{p_{\text{data}}(\mathbf{x})}[q_\phi(\mathbf{z}|\mathbf{x})]$ . A mis-  
 145 match could create “holes” [76, 41, 3] in the prior that the aggregate posterior fails to cover during  
 146 training, resulting in worse sample quality, as many latent variables used during generation are likely  
 147 to never have been trained on. We formalize this notion in the following definition.

148 **Definition 1** (Prior hole). *Let  $p(\mathbf{z}), q(\mathbf{z})$  be two distributions with  $\text{supp}(q) \subseteq \text{supp}(p)$ . We say that*  
 149  *$q$  has an  $(\epsilon, \delta)$ -prior hole with respect to (the prior)  $p$  for  $\epsilon, \delta \in (0, 1)$ ,  $\delta > \epsilon$ , if there exists a set*  
 150  *$S \in \text{supp}(P)$ , such that  $\int_S p(\mathbf{z})d\mathbf{z} \geq \delta$  and  $\int_S q(\mathbf{z})d\mathbf{z} \leq \epsilon$ .*

151 Intuitively, if  $q_\phi(\mathbf{z})$  has a prior hole with large  $\delta$  and small  $\epsilon$  (e.g., inversely proportional to the number  
 152 of training samples), then it is very likely that latent variables within the hole are never seen during  
 153 training (small  $\epsilon$ ), yet frequently used to produce samples (large  $\delta$ ). Most existing methods address  
 154 this problem by optimizing certain statistical divergences between  $q_\phi(\mathbf{z})$  and  $p_\theta(\mathbf{z})$ , such as the KL  
 155 divergence or Wasserstein distance [84]. However, we argue in the following statement that prior  
 156 holes might not be eliminated even if we optimize certain divergence values to be reasonably low,  
 157 especially when  $q_\phi(\mathbf{z})$  is very flexible. We present the formal statement and proof in Appendix B.2.  
 158

159 **Theorem 2.** (informal) *Let  $p_\theta(\mathbf{z}) = \mathcal{N}(0, 1)$ . For any  $\epsilon > 0$ ,*  
 160 *there exists a distribution  $q_\phi(\mathbf{z})$  with an  $(\epsilon, 0.49)$ -prior hole, such that*  
 161  *$D_{\text{KL}}(q_\phi||p_\theta) \leq \log 2^2$  and  $W_2(q_\phi, p_\theta) < \gamma$  for any  $\gamma > 0$ , where  $W_2$  is*  
 162 *the 2-Wasserstein distance.*



163 *Proof.* (sketch) We construct a  $q_\phi$  that satisfies these properties (top-right figure). First, we truncate  
 164 the Gaussian and divide them into regions with same probability mass; then we support  $q_\phi$  over half  
 165 of these regions (so  $\delta > 0.49$ ); finally, we show that the divergences are small enough.  $\square$

166 In contrast to addressing prior holes by optimization, diffusion models eliminate prior holes by  
 167 construction, since the diffusion process from  $\mathbf{z}^{(1)}$  to  $\mathbf{z}^{(0)}$  is constructed such that the distribution  
 168 of  $\mathbf{z}^{(\alpha)}$  always converges to a standard Gaussian as  $\alpha \rightarrow 0$ . As a result, the distribution of latent  
 169 variables used during training is arbitrarily close to that used in generation<sup>3</sup>, which is also the case in  
 170 GANs. Therefore, our argument provides an explanation as to why we observe better sample quality  
 171 results from GANs and diffusion models than VAEs and NFs.

172 **5 Few-shot Conditional Generation with D2C**

173 In this section, we discuss how D2C can be used to learn to perform conditional generation from  
 174 few-shot supervision. We note that D2C is only trained on images and not with any other data  
 175 modalities (e.g., image-text pairs [71]) or supervision techniques (e.g., meta-learning [20, 6]).

<sup>2</sup>This is reasonably low for realistic VAE models (NVAE [87] reports a KL divergence of around 2810 nats).  
<sup>3</sup>We expand this argument in Appendix B.2.

176 **Algorithm** We describe the general algorithm  
 177 for conditional generation from a few images  
 178 in Algorithm 1, and detailed implementations  
 179 in Appendix C. With a model over the latent  
 180 space (denoted as  $r_\psi(\mathbf{c}|\mathbf{z}^{(1)})$ ), we draw condi-  
 181 tional latents from an unnormalized distribution  
 182 with the diffusion prior (line 4). This can be  
 183 implemented in many ways such as rejection  
 184 sampling or Langevin dynamics [63, 82, 25].

---

**Algorithm 1** Conditional generation with D2C

---

- 1: **Input**  $n$  examples  $\{(\mathbf{x}_i, \mathbf{c}_i)\}_{i=1}^n$ , property  $\mathbf{c}$ .
  - 2: Acquire latents  $\mathbf{z}_i^{(1)} \sim q_\phi(\mathbf{z}^{(1)}|\mathbf{x})$  for  $i \in [n]$ ;
  - 3: Train model  $r_\psi(\mathbf{c}|\mathbf{z}^{(1)})$  over  $\{(\mathbf{z}_i^{(1)}, \mathbf{c}_i)\}_{i=1}^n$
  - 4: Sample latents with  $\hat{\mathbf{z}}^{(1)} \sim r_\psi(\mathbf{c}|\mathbf{z}^{(1)}) \cdot p_\theta^{(1)}(\mathbf{z}^{(1)})$   
(unnormalized);
  - 5: Decode  $\hat{\mathbf{x}} \sim p_\theta(\mathbf{x}|\hat{\mathbf{z}}^{(1)})$ .
  - 6: **Output**  $\hat{\mathbf{x}}$ .
- 

185 **Conditions from labeled examples** Given a few labeled examples, we wish to produce diverse  
 186 samples with a certain label. For labeled examples we can directly train a classifier over the  
 187 latent space, which we denote as  $r_\psi(\mathbf{c}|\mathbf{z}^{(1)})$  with  $\mathbf{c}$  being the class label and  $\mathbf{z}^{(1)}$  being the  
 188 representation of  $\mathbf{x}$  from  $q_\phi(\mathbf{z}^{(1)}|\mathbf{x})$ . If these examples do not have labels (*i.e.*, we merely want to  
 189 generate new samples similar to given ones), we can train a positive-unlabeled (PU) classifier [31]  
 190 where we assign “positive” to the new examples and “unlabeled” to training data. Then we use  
 191 the classifier with the diffusion model  $p_\theta(\mathbf{z}^{(1)}|\mathbf{z}^{(0)})$  to produce suitable values of  $\mathbf{z}^{(1)}$ , such as by  
 192 rejecting samples from the diffusion model that has a small  $r_\psi(\mathbf{c}|\mathbf{z}^{(1)})$ .

193 **Conditions from manipulation constraints** Given a few labeled examples, here we wish to learn  
 194 how to manipulate images. Specifically, we condition over the event that “ $\mathbf{x}$  has label  $\mathbf{c}$  but is  
 195 similar to image  $\bar{\mathbf{x}}$ ”. Here  $r_\psi(\mathbf{c}|\mathbf{z}^{(1)})$  is the unnormalized product between the classifier conditional  
 196 probability and closeness to the latent  $\bar{\mathbf{z}}^{(1)}$  of  $\bar{\mathbf{x}}$  (*e.g.*, measured with RBF kernel). We implement line  
 197 4 of Alg. 1 with a Lanvegin-like procedure where we take a gradient step with respect to the classifier  
 198 probability and then correct this gradient step with the diffusion model. Unlike many GAN-based  
 199 methods [12, 69, 92, 43, 93], D2C does not need to optimize an inversion procedure at evaluation  
 200 time, and thus the latent value is much faster to compute; D2C is also better at retaining fine-grained  
 201 features of the original image due to the reconstruction loss.

## 202 6 Related Work

203 **Latent variable generative models** Most deep generative models explicitly define a latent rep-  
 204 resentation, except for some energy-based models [39, 29] and autoregressive models [89, 88, 10].  
 205 Unlike VAEs and NFs, GANs do not explicitly define an inference model and instead optimize a  
 206 two-player game. In terms of sample quality, GANs currently achieve superior performance over  
 207 VAEs and NFs, but they can be difficult to invert even with additional optimization [45, 95, 7]. This  
 208 can be partially addressed by training reconstruction-based losses with GANs [51, 52]. Moreover,  
 209 the GAN training procedure can be unstable [9, 8, 60], lack a informative objective for measuring  
 210 progress [4], and struggle with discrete data [97]. Diffusion models [25] achieves high sample quality  
 211 without adversarial training, but its latent dimension must be equal to the image dimension.

212 **Addressing posterior mismatch in VAEs** Most methods address this mismatch problem by im-  
 213 proving inference models [61, 48, 85], prior models [86, 3, 83], or objective functions [98, 99, 100,  
 214 1, 56]; all these approaches optimize the posterior model to be close to the prior. In Section 4.2,  
 215 we explain why these approaches do not necessarily remove large “prior holes”, so their sample  
 216 qualities remain relatively poor even after many layers [87, 18]. Other methods adopt a “two-stage”  
 217 approach [21], which fits a generative model over the latent space of autoencoders [91, 72, 24, 71].

218 **Conditional generation with unconditional models** To perform conditional generation over an  
 219 unconditional LVGM, most methods assume access to a discriminative model (*e.g.*, a classifier); the  
 220 latent space of the generator is then modified to change the outputs of the discriminative model. The  
 221 discriminative model can operate on either the image space [63, 67, 25] or the latent space [77, 93]. For  
 222 image space discriminative models, plug-and-play generative networks [63] control the attributes of  
 223 generated images via Langevin dynamics [75]; these ideas are also explored in diffusion models [82].  
 224 Image manipulation methods are based on GANs often operate with latent space discriminators [77,  
 225 93]. However, these methods have some trouble manipulating real images because of imperfect  
 226 reconstruction [102, 7]. This is not a problem in D2C since a reconstruction objective is optimized.



Figure 3: Generated samples on CIFAR-10 (left), fMoW (mid), and FFHQ  $256 \times 256$  (right).

## 227 7 Experiments

228 We examine the conditional and unconditional generation qualities of D2C over CIFAR-10 [50],  
 229 CIFAR-100 [50], fMoW [19], CelebA-64 [55], CelebA-HQ-256 [45], and FFHQ-256 [46]. Our D2C  
 230 implementation is based on the state-of-the-art NVAE [87] autoencoder structure, the U-Net diffusion  
 231 model [40], and the MoCo-v2 contrastive representation learning method [15]. We keep the diffusion  
 232 series hyperparameter  $\{\alpha_i\}_{i=1}^T$  identical to ensure a fair comparison with different diffusion models.  
 233 For the contrastive weight hyperparameter  $\lambda$  in Equation (4), we consider the value of  $\lambda = 10^{-4}$   
 234 based on the relative scale between the  $L_C$  and  $L_{D2}$ ; we find that the results are relatively insensitive  
 235 to  $\lambda$ . We use 100 diffusion steps for DDIM and D2C unless mentioned otherwise, as running with  
 236 longer steps is not computationally economical despite tiny gains in FID [80]. We include additional  
 237 training details, such as architectures, optimizers and learning rates in Appendix C.

Table 2: Quality of representations and generations with LVGMs.

Model	CIFAR-10			CIFAR-100			fMoW		
	FID ↓	MSE ↓	Acc ↑	FID ↓	MSE ↓	Acc ↑	FID ↓	MSE ↓	Acc ↑
NVAE [87]	36.4	<b>0.25</b>	18.8	42.5	0.53	4.1	82.25	<b>0.30</b>	27.7
DDIM [80]	<b>4.16</b>	2.5	22.5	<b>10.16</b>	3.2	2.2	<b>37.74</b>	3.0	23.5
D2C (Ours)	10.15	0.76	<b>76.02</b>	14.62	<b>0.44</b>	<b>42.75</b>	44.7	2.33	<b>66.9</b>

### 238 7.1 Unconditional generation

239 For unconditional generation, we measure the sample quality of images using the Fréchet Inception  
 240 Distance (FID, [38]) with 50,000 images. In particular, we extensively evaluate NVAE [87] and  
 241 DDIM [80], a competitive VAE model and a competitive diffusion model as baselines because we  
 242 can directly obtain features from them without additional optimization steps<sup>4</sup>. For them, we report  
 243 mean-squared reconstruction error (MSE, summed over all pixels, pixels normalized to  $[0, 1]$ ) and  
 244 linear classification accuracy (Acc., measured in percentage) over  $\mathbf{z}_1$  features for the test set.

245 We report sample quality results<sup>5</sup> in Tables 2, and 3. For FID, we outperform NVAE in all datasets and  
 246 outperform DDIM on CelebA-64 and CelebA-HQ-256, which suggests our results are competitive  
 247 with state-of-the-art non-adversarial generative models. In Table 2, we additionally compare NVAE,  
 248 DDIM and D2C in terms of reconstruction and linear classification accuracy. As all three methods  
 249 contain reconstruction losses, the MSE values are low and comparable. However, D2C enjoys much  
 250 better linear classification accuracy than the other two thanks to the contrastive SSL component. We  
 251 further note that training the same contrastive SSL method without  $L_{D2}$  achieves slightly higher  
 252 78.3% accuracy on CIFAR-10. We tried improving this via ResNet [36] encoders, but this significantly  
 253 increased reconstruction error, possibly due to loss of information in average pooling layers.

<sup>4</sup>For DDIM, the latent representations  $\mathbf{x}^{(0)}$  are obtained by reversing the neural ODE process.

<sup>5</sup>Due to space limits, we place additional CIFAR-10 results in Appendix D.

Table 3: FID scores over different faces dataset with LVGMs.

Model	CelebA-64	CelebA-HQ-256	FFHQ-256
NVAE [87]	13.48	40.26	26.02
DDIM [80]	6.53	25.6	-
D2C (Ours)	<b>5.7</b>	<b>18.74</b>	<b>13.04</b>

Table 4: Sample quality as a function of diffusion steps.

Steps	CIFAR-10			CIFAR-100			CelebA-64		
	10	50	100	10	50	100	10	50	100
DDPM [40]	41.07	8.01	5.78	50.27	21.37	16.72	33.12	18.48	13.93
DDIM [80]	<b>13.36</b>	<b>4.67</b>	<b>4.16</b>	<b>23.34</b>	<b>11.69</b>	<b>10.16</b>	17.33	9.17	6.53
D2C (Ours)	17.71	10.11	10.15	23.16	14.62	14.46	<b>17.32</b>	<b>6.8</b>	<b>5.7</b>

## 254 7.2 Few-shot conditional generation from examples

255 We demonstrate the advantage of D2C representations by performing few-shot conditional generation  
 256 over labels. We consider two types of labeled examples: one has binary labels for which we train  
 257 a binary classifier; the other is positive-only labeled (*e.g.*, images of female class) for which we  
 258 train a PU classifier. Our goal here is to generate a diverse group of images with a certain label. We  
 259 evaluate and compare three models: D2C, NVAE and DDIM. We train a classifier  $r_\psi(c|z)$  over the  
 260 latent space of these models; we also train a image space classifier and use it with DDIM (denoted as  
 261 DDIM-I). We run Algorithm 1 for these models, where line 4 is implemented via rejection sampling.  
 262 As our goal is to compare different models, we leave more sophisticated methods [25] as future work.

263 We consider performing 8 conditional generation tasks over CelebA-64 with 2 binary classifiers  
 264 (trained over 100 samples, 50 for each class) and 4 PU classifiers (trained over 100 positively labeled  
 265 and 10k unlabeled samples). We also report a “naive” approach where we use all the training  
 266 images (regardless of labels) and compute its FID with the corresponding subset of images (*e.g.*, all  
 267 images versus blond images). In Table 5, we report the FID score between generated images (5k  
 268 samples) and real images of the corresponding label. These results suggest that D2C outperforms the  
 269 other approaches, and is the only one that performs better than the “naive” approach in most cases,  
 illustrating the advantage of contrastive representations for few-shot conditional generation.

Table 5: FID scores for few-shot conditional generation with various types of labeled examples. Naive performs very well for non-blond due to class percentages.

Method	Classes (% in train set)	D2C	DDIM	NVAE	DDIM-I	Naive
Binary	Male (42%)	<b>13.44</b>	38.38	41.07	29.03	26.34
	Female (58%)	<b>9.51</b>	19.25	16.57	15.17	18.72
	Blond (15%)	<b>17.61</b>	31.39	31.24	29.09	27.51
	Non-Blond (85%)	<b>8.94</b>	9.67	16.73	19.76	3.77
PU	Male (42%)	<b>16.39</b>	37.03	42.78	19.60	26.34
	Female (58%)	<b>12.21</b>	15.42	18.36	14.96	18.72
	Blond (15%)	<b>10.09</b>	30.20	31.06	76.52	27.51
	Non-Blond (85%)	<b>9.09</b>	9.70	17.98	9.90	3.77

270

## 271 7.3 Few-shot conditional generation from manipulation constraints

272 Finally, we consider image manipulation where we use binary classifiers that are learned over 50  
 273 labeled instances for each class. We perform Amazon Mechanical Turk (AMT) evaluations over  
 274 two attributes in the CelebA-256 dataset, *blond* and *red lipstick*, over D2C, DDIM, NVAE and



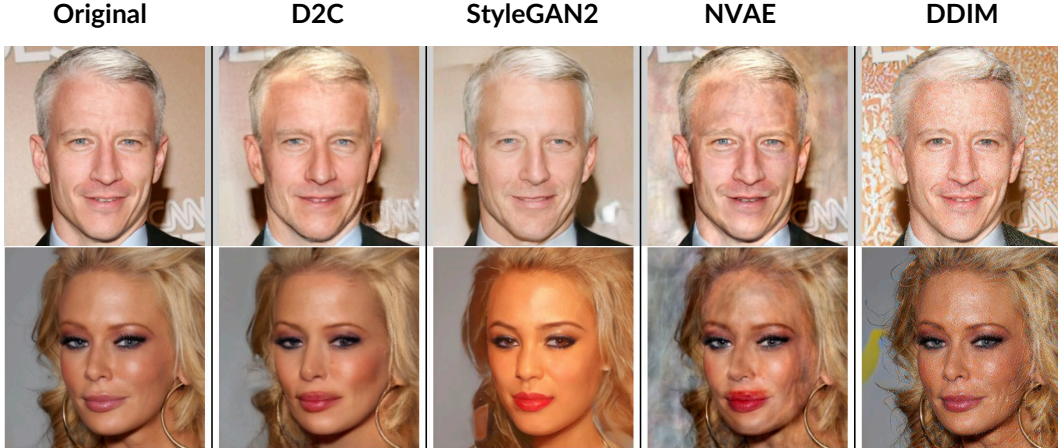


Figure 4: Image manipulation results for *blond* (top) and *red lipstick* (bottom). D2C is better than StyleGAN2 at preserving details of the original image, such as eyes, earrings, and background.

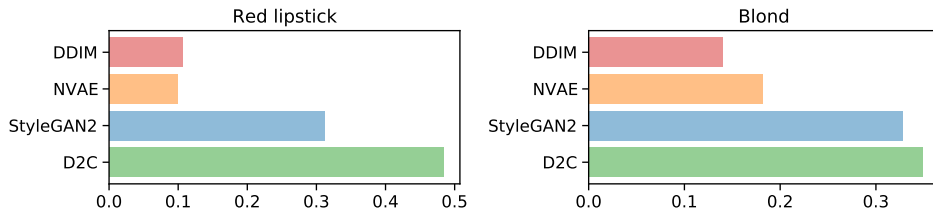


Figure 5: AMT evaluation over image manipulations.  $x$ -axis shows the percentage that the evaluator selects the image generated from the corresponding model out of 4 images from each model.

275 StyleGAN2 [47] (see Figure 4). The evaluation is double-blinded: neither we nor the evaluators know  
 276 the correspondence between generated image and underlying model during the study. We include  
 277 more details (algorithm, setup and human evaluation) in Appendix C and additional qualitative results  
 278 (such as *beard* and *gender* attributes) in Appendix D.

279 In Figure 5, we show the percentage of manipulations preferred by AMT evaluators for each model;  
 280 D2C slightly outperforms StyleGAN2 for *blond* and significantly outperforms StyleGAN2 for *red*  
 281 *lipstick*. When we compare D2C with only StyleGAN2, D2C is preferred over 51.5% for *blond*  
 282 and 60.8% for *red lipstick*. An additional advantage of D2C is that the manipulation is much faster  
 283 than StyleGAN2, since the latter requires additional optimization over the latent space to improve  
 284 reconstruction [101]. On the same Nvidia 1080Ti GPU, it takes 0.013 seconds to obtain the latent  
 285 code in D2C, while the same takes 8 seconds [101] for StyleGAN2 (615× slower). As decoding is  
 286 very fast for both models, D2C generations are around two orders of magnitude faster to produce.

## 287 8 Discussions and Limitations

288 We introduced D2C, a VAE-based generative model with a latent space suitable for few-shot condi-  
 289 tional generation. To our best knowledge, our model is the first unconditional VAE to demonstrate  
 290 superior image manipulation performance than StyleGAN2, which is surprising given our use of a  
 291 regular NVAE architecture. We believe that with better architectures, such as designs from Style-  
 292 GAN2 or Transformers [42], D2C can achieve even better performance. It is also interesting to  
 293 formally investigate the integration between D2C and other types of conditions on the latent space, as  
 294 well as training D2C in conjunction with other domains and data modalities, such as text [71], in a  
 295 fashion that is similar to semi-supervised learning. Nevertheless, we note that our model have to be  
 296 used properly in order to mitigate potential negative societal impacts, such as deep fakes.

297 **References**

- 298 [1] Alexander A Alemi, Ben Poole, Ian Fischer, Joshua V Dillon, Rif A Saurous, and Kevin  
299 Murphy. Fixing a broken ELBO. *arXiv preprint arXiv:1711.00464*, November 2017.
- 300 [2] Jyoti Aneja, Harsh Agrawal, Dhruv Batra, and Alexander Schwing. Sequential latent spaces  
301 for modeling the intention during diverse image captioning. In *Proceedings of the IEEE/CVF*  
302 *International Conference on Computer Vision*, pages 4261–4270, 2019.
- 303 [3] Jyoti Aneja, Alexander Schwing, Jan Kautz, and Arash Vahdat. NCP-VAE: Variational  
304 autoencoders with noise contrastive priors. *arXiv preprint arXiv:2010.02917*, October 2020.
- 305 [4] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *arXiv preprint*  
306 *arXiv:1701.07875*, January 2017.
- 307 [5] Muhammad Asim, Ali Ahmed, and Paul Hand. Invertible generative models for inverse  
308 problems: mitigating representation error and dataset bias. *arXiv preprint arXiv:1905.11672*,  
309 May 2019.
- 310 [6] Sergey Bartunov and Dmitry Vetrov. Few-shot generative modelling with generative matching  
311 networks. In Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First*  
312 *International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of*  
313 *Machine Learning Research*, pages 670–678. PMLR, 2018.
- 314 [7] David Bau, Jun-Yan Zhu, Jonas Wulff, William Peebles, Hendrik Strobelt, Bolei Zhou, and  
315 Antonio Torralba. Seeing what a gan cannot generate. In *Proceedings of the IEEE/CVF*  
316 *International Conference on Computer Vision*, pages 4502–4511, 2019.
- 317 [8] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity  
318 natural image synthesis. *arXiv preprint arXiv:1809.11096*, September 2018.
- 319 [9] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with  
320 introspective adversarial networks. *arXiv preprint arXiv:1609.07093*, 2016.
- 321 [10] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-  
322 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,  
323 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M  
324 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz  
325 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec  
326 Radford, Ilya Sutskever, and Dario Amodei. Language models are Few-Shot learners. *arXiv*  
327 *preprint arXiv:2005.14165*, May 2020.
- 328 [11] Emmanuel J Candes, Justin K Romberg, and Terence Tao. Stable signal recovery from  
329 incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics:*  
330 *A Journal Issued by the Courant Institute of Mathematical Sciences*, 59(8):1207–1223, 2006.
- 331 [12] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement  
332 networks. In *ICCV*, 2017.
- 333 [13] Ricky T Q Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary  
334 differential equations. *arXiv preprint arXiv:1806.07366*, June 2018.
- 335 [14] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
336 for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, February  
337 2020.
- 338 [15] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum  
339 contrastive learning. *arXiv preprint arXiv:2003.04297*, March 2020.
- 340 [16] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. *arXiv*  
341 *preprint arXiv:2011.10566*, November 2020.
- 342 [17] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training Self-Supervised  
343 vision transformers. *arXiv preprint arXiv:2104.02057*, April 2021.

- 344 [18] Rewon Child. Very deep vaes generalize autoregressive models and can outperform them on  
345 images. *arXiv preprint arXiv:2011.10650*, 2020.
- 346 [19] Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the  
347 world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
348 pages 6172–6180, 2018.
- 349 [20] Louis Clouâtre and Marc Demers. Figr: Few-shot image generation with reptile. *arXiv preprint*  
350 *arXiv:1901.02199*, 2019.
- 351 [21] Bin Dai and David Wipf. Diagnosing and enhancing VAE models. *arXiv preprint*  
352 *arXiv:1903.05789*, March 2019.
- 353 [22] Giannis Daras, Joseph Dean, Ajil Jalal, and Alexandros G Dimakis. Intermediate layer opti-  
354 mization for inverse problems using deep generative models. *arXiv preprint arXiv:2102.07364*,  
355 February 2021.
- 356 [23] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
357 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,  
358 October 2018.
- 359 [24] Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya  
360 Sutskever. Jukebox: A generative model for music. *arXiv preprint arXiv:2005.00341*, 2020.
- 361 [25] Prafulla Dhariwal and Alex Nichol. Diffusion models beat GANs on image synthesis. *arXiv*  
362 *preprint arXiv:2105.05233*, May 2021.
- 363 [26] L Dinh, D Krueger, and Y Bengio. NICE: Non-linear independent components estimation.  
364 *arXiv preprint arXiv:1410.8516*, 2014.
- 365 [27] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP.  
366 *arXiv preprint arXiv:1605.08803*, May 2016.
- 367 [28] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv*  
368 *preprint arXiv:1605.09782*, May 2016.
- 369 [29] Yilun Du and Igor Mordatch. Implicit generation and generalization in Energy-Based models.  
370 *arXiv preprint arXiv:1903.08689*, March 2019.
- 371 [30] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Mar-  
372 tin Arjovsky, and Aaron Courville. Adversarially learned inference. *arXiv preprint*  
373 *arXiv:1606.00704*, June 2016.
- 374 [31] Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In  
375 *14th ACM SIGKDD*, pages 213–220, 2008.
- 376 [32] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
377 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z Ghahramani,  
378 M Welling, C Cortes, N D Lawrence, and K Q Weinberger, editors, *Advances in Neural*  
379 *Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.
- 380 [33] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena  
381 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Ghesh-  
382 laghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your  
383 own latent: A new approach to Self-Supervised learning. *arXiv preprint arXiv:2006.07733*,  
384 June 2020.
- 385 [34] Aditya Grover, Aaron Zweig, and Stefano Ermon. Graphite: Iterative generative modeling of  
386 graphs. In *International conference on machine learning*, pages 2434–2444. PMLR, 2019.
- 387 [35] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast  
388 for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, November  
389 2019.

- 390 [36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
391 recognition. *arXiv preprint arXiv:1512.03385*, December 2015.
- 392 [37] Olivier Henaff. Data-efficient image recognition with contrastive predictive coding. In  
393 *International Conference on Machine Learning*, pages 4182–4192. PMLR, 2020.
- 394 [38] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
395 GANs trained by a two Time-Scale update rule converge to a local nash equilibrium. *arXiv*  
396 *preprint arXiv:1706.08500*, June 2017.
- 397 [39] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural*  
398 *computation*, 14(8):1771–1800, August 2002.
- 399 [40] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv*  
400 *preprint arXiv:2006.11239*, June 2020.
- 401 [41] Matthew D Hoffman and Matthew J Johnson. Elbo surgery: yet another way to carve up  
402 the variational evidence lower bound. In *Workshop in Advances in Approximate Bayesian*  
403 *Inference, NIPS*, volume 1, page 2. approximateinference.org, 2016.
- 404 [42] Drew A Hudson and C Lawrence Zitnick. Generative adversarial transformers. *arXiv preprint*  
405 *arXiv:2103.01209*, 2021.
- 406 [43] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation  
407 with conditional adversarial networks. In *CVPR*, 2017.
- 408 [44] Zahra Kadkhodaie and Eero P Simoncelli. Solving linear inverse problems using the prior  
409 implicit in a denoiser. *arXiv preprint arXiv:2007.13640*, July 2020.
- 410 [45] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs  
411 for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, October 2017.
- 412 [46] Tero Karras, Samuli Laine, and Timo Aila. A Style-Based generator architecture for generative  
413 adversarial networks. *arXiv preprint arXiv:1812.04948*, December 2018.
- 414 [47] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila.  
415 Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF*  
416 *Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020.
- 417 [48] Diederik P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max  
418 Welling. Improved variational inference with inverse autoregressive flow. In D D Lee,  
419 M Sugiyama, U V Luxburg, I Guyon, and R Garnett, editors, *Advances in Neural Information*  
420 *Processing Systems 29*, pages 4743–4751. Curran Associates, Inc., 2016.
- 421 [49] Diederik P Kingma and Max Welling. Auto-Encoding variational bayes. *arXiv preprint*  
422 *arXiv:1312.6114v10*, December 2013.
- 423 [50] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep  
424 convolutional neural networks. In F Pereira, C J C Burges, L Bottou, and K Q Weinberger,  
425 editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran  
426 Associates, Inc., 2012.
- 427 [51] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther.  
428 Autoencoding beyond pixels using a learned similarity metric. In *International conference on*  
429 *machine learning*, pages 1558–1566. PMLR, 2016.
- 430 [52] C Li, H Liu, C Chen, Y Pu, L Chen, and others. Alice: Towards understanding adversarial  
431 learning for joint distribution matching. *Advances in neural information processing systems*,  
432 2017.
- 433 [53] Chunyuan Li, Xiang Gao, Yuan Li, Baolin Peng, Xiujun Li, Yizhe Zhang, and Jianfeng Gao.  
434 Optimus: Organizing sentences via pre-trained modeling of a latent space. *arXiv preprint*  
435 *arXiv:2004.04092*, 2020.

- 436 [54] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,  
437 Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In  
438 *European conference on computer vision*, pages 740–755. Springer, 2014.
- 439 [55] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in  
440 the wild. In *Proceedings of the IEEE international conference on computer vision*, pages  
441 3730–3738, 2015.
- 442 [56] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey.  
443 Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.
- 444 [57] Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. Generating  
445 images from captions with attention. *arXiv preprint arXiv:1511.02793*, 2015.
- 446 [58] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint*  
447 *arXiv:1411.1784*, November 2014.
- 448 [59] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation  
449 with diffusion models. *arXiv preprint arXiv:2103.16091*, March 2021.
- 450 [60] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normaliza-  
451 tion for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, February 2018.
- 452 [61] Shakir Mohamed and Balaji Lakshminarayanan. Learning in implicit generative models. *arXiv*  
453 *preprint arXiv:1610.03483*, October 2016.
- 454 [62] Alex Najibi. Racial Discrimination in Face Recognition Technology, 2020.
- 455 [63] Anh Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, and Jason Yosinski. Plug & play  
456 generative networks: Conditional iterative generation of images in latent space. In *Proceedings*  
457 *of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4467–4477, 2017.
- 458 [64] Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon.  
459 Permutation invariant graph generation via score-based generative modeling. In *International*  
460 *Conference on Artificial Intelligence and Statistics*, pages 4474–4484. PMLR, 2020.
- 461 [65] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving  
462 jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016.
- 463 [66] Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei A Efros,  
464 and Richard Zhang. Swapping autoencoder for deep image manipulation. *arXiv preprint*  
465 *arXiv:2007.00653*, July 2020.
- 466 [67] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. StyleCLIP:  
467 Text-Driven manipulation of StyleGAN imagery. *arXiv preprint arXiv:2103.17249*, March  
468 2021.
- 469 [68] Wei Ping, Kainan Peng, Kexin Zhao, and Zhao Song. Waveflow: A compact flow-based model  
470 for raw audio. In *International Conference on Machine Learning*, pages 7706–7716. PMLR,  
471 2020.
- 472 [69] Tiziano Portenier, Qiyang Hu, Attila Szabó, Siavash Arjomand Bigdeli, Paolo Favaro, and  
473 Matthias Zwicker. Faceshop: Deep sketch-based face image editing. *ACM Transactions on*  
474 *Graphics*, 37(4), 2018.
- 475 [70] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini  
476 Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and  
477 Ilya Sutskever. Learning transferable visual models from natural language supervision. *arXiv*  
478 *preprint arXiv:2103.00020*, February 2021.
- 479 [71] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford,  
480 Mark Chen, and Ilya Sutskever. Zero-Shot Text-to-Image generation. *arXiv preprint*  
481 *arXiv:2102.12092*, February 2021.

- 482 [72] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images  
483 with vq-vae-2. *arXiv preprint arXiv:1906.00446*, 2019.
- 484 [73] Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak  
485 Lee. Generative adversarial text to image synthesis. In *International Conference on Machine*  
486 *Learning*, pages 1060–1069. PMLR, 2016.
- 487 [74] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows.  
488 *arXiv preprint arXiv:1505.05770*, May 2015.
- 489 [75] Gareth O Roberts and Jeffrey S Rosenthal. Optimal scaling of discrete approximations to  
490 langevin diffusions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,  
491 60(1):255–268, 1998.
- 492 [76] Mihaela Rosca, Balaji Lakshminarayanan, and Shakir Mohamed. Distribution matching in  
493 variational inference. *arXiv preprint arXiv:1802.06847*, February 2018.
- 494 [77] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for  
495 semantic face editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
496 *Pattern Recognition*, pages 9243–9252, 2020.
- 497 [78] Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsu-  
498 pervised learning using nonequilibrium thermodynamics. *arXiv preprint arXiv:1503.03585*,  
499 March 2015.
- 500 [79] Jiaming Song and Stefano Ermon. Multi-label contrastive predictive coding. *arXiv preprint*  
501 *arXiv:2007.09852*, 2020.
- 502 [80] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
503 *preprint arXiv:2010.02502*, 2020.
- 504 [81] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data  
505 distribution. *arXiv preprint arXiv:1907.05600*, July 2019.
- 506 [82] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and  
507 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv*  
508 *preprint arXiv:2011.13456*, 2020.
- 509 [83] Hiroshi Takahashi, Tomoharu Iwata, Yuki Yamanaka, Masanori Yamada, and Satoshi Yagi.  
510 Variational autoencoder with implicit optimal priors. *Proceedings of the AAAI Conference on*  
511 *Artificial Intelligence*, 33(01):5066–5073, July 2019.
- 512 [84] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein  
513 Auto-Encoders. *arXiv preprint arXiv:1711.01558*, November 2017.
- 514 [85] Jakub M Tomczak and Max Welling. Improving variational auto-encoders using householder  
515 flow. *arXiv preprint arXiv:1611.09630*, 2016.
- 516 [86] Jakub M Tomczak and Max Welling. VAE with a VampPrior. *arXiv preprint arXiv:1705.07120*,  
517 May 2017.
- 518 [87] Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. *arXiv*  
519 *preprint arXiv:2007.03898*, July 2020.
- 520 [88] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex  
521 Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. WaveNet: A generative  
522 model for raw audio. *arXiv preprint arXiv:1609.03499*, September 2016.
- 523 [89] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural  
524 networks. *arXiv preprint arXiv:1601.06759*, January 2016.
- 525 [90] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive  
526 predictive coding. *arXiv preprint arXiv:1807.03748*, July 2018.

- 527 [91] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation  
528 learning. *arXiv preprint arXiv:1711.00937*, November 2017.
- 529 [92] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro.  
530 High-resolution image synthesis and semantic manipulation with conditional gans. In *CVPR*,  
531 2018.
- 532 [93] Weihao Xia, Yujia Yang, Jing-Hao Xue, and Baoyuan Wu. TediGAN: Text-Guided diverse  
533 face image generation and manipulation. 2021.
- 534 [94] Zhe Xie, Chengxuan Liu, Yichi Zhang, Hongtao Lu, Dong Wang, and Yue Ding. Adversarial  
535 and contrastive variational autoencoder for sequential recommendation. *arXiv preprint*  
536 *arXiv:2103.10693*, March 2021.
- 537 [95] Yinghao Xu, Yujun Shen, Jiapeng Zhu, Ceyuan Yang, and Bolei Zhou. Generative hierarchical  
538 features from synthesizing images. *arXiv e-prints*, pages arXiv–2007, 2020.
- 539 [96] Jiakuan You, Bowen Liu, Rex Ying, Vijay Pande, and Jure Leskovec. Graph convolutional pol-  
540 icy network for Goal-Directed molecular graph generation. *arXiv preprint arXiv:1806.02473*,  
541 June 2018.
- 542 [97] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: Sequence generative adversarial  
543 nets with policy gradient. In *Proceedings of the AAAI conference on artificial intelligence*,  
544 volume 31, 2017.
- 545 [98] Shengjia Zhao, Jiaming Song, and Stefano Ermon. InfoVAE: Information maximizing varia-  
546 tional autoencoders. *arXiv preprint arXiv:1706.02262*, June 2017.
- 547 [99] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Towards deeper understanding of variational  
548 autoencoding models. *arXiv preprint arXiv:1702.08658*, February 2017.
- 549 [100] Shengjia Zhao, Jiaming Song, and Stefano Ermon. A lagrangian perspective on latent variable  
550 generative models. In *Proc. 34th Conference on Uncertainty in Artificial Intelligence*, 2018.
- 551 [101] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain gan inversion for real image  
552 editing. In *European Conference on Computer Vision*, pages 592–608. Springer, 2020.
- 553 [102] Jiapeng Zhu, Deli Zhao, Bolei Zhou, and Bo Zhang. Lia: Latently invertible autoencoder with  
554 adversarial learning. 2019.

555 **Checklist**

- 556 1. For all authors...
- 557 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
558 contributions and scope? [Yes]
- 559 (b) Did you describe the limitations of your work? [Yes] In Section 8.
- 560 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Details in  
561 Section 8 and Appendix E.
- 562 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
563 them? [Yes]
- 564 2. If you are including theoretical results...
- 565 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 566 (b) Did you include complete proofs of all theoretical results? [Yes] In the Appendix.
- 567 3. If you ran experiments...
- 568 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
569 mental results (either in the supplemental material or as a URL)? [Yes]
- 570 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
571 were chosen)? [No] We report everything except data splits for standard datasets such  
572 as CIFAR-10.
- 573 (c) Did you report error bars (e.g., with respect to the random seed after running ex-  
574 periments multiple times)? [No] The experiments are simply too expensive (and  
575 environmentally harmful) to be run for multiple times in order to evaluate error bars;  
576 this is also common practice in the generative modeling literature.
- 577 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
578 of GPUs, internal cluster, or cloud provider)? [Yes] In Appendix C.
- 579 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 580 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 581 (b) Did you mention the license of the assets? [N/A] CIFAR-10 does not have a license.
- 582 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
583 No new assets are included.
- 584 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
585 using/curating? [N/A] All the datasets are public.
- 586 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
587 information or offensive content? [N/A] All the datasets are public.
- 588 5. If you used crowdsourcing or conducted research with human subjects...
- 589 (a) Did you include the full text of instructions given to participants and screenshots, if  
590 applicable? [Yes] In Appendix C.
- 591 (b) Did you describe any potential participant risks, with links to Institutional Review  
592 Board (IRB) approvals, if applicable? [N/A] No biomedical data is used.
- 593 (c) Did you include the estimated hourly wage paid to participants and the total amount  
594 spent on participant compensation? [Yes] In Appendix C.