
Generative Visual Prompt: Unifying Distributional Control of Pre-Trained Generative Models

Chen Henry Wu, Saman Motamed, Shaunak Srivastava, Fernando De la Torre
Carnegie Mellon University
henrychenwu@cmu.edu

Abstract

Generative models (e.g., GANs and diffusion models) learn the underlying data distribution in an unsupervised manner. However, many applications of interest require sampling from a specific region of the generative model’s output space or evenly over a range of characteristics. To allow efficient sampling in these scenarios, we propose Generative Visual Prompt (PromptGen), a framework for distributional control over pre-trained generative models by incorporating knowledge of arbitrary off-the-shelf models. PromptGen defines control as an energy-based model (EBM) and samples images in a feed-forward manner by approximating the EBM with invertible neural networks, avoiding optimization at inference. We demonstrate how PromptGen can control several generative models (e.g., StyleGAN2, StyleNeRF, diffusion autoencoder, and NVAE) using various off-the-shelf models: (1) with the CLIP model, PromptGen can sample images guided by text, (2) with image classifiers, PromptGen can de-bias generative models across a set of attributes, and (3) with inverse graphics models, PromptGen can sample images of the same identity in different poses. (4) Finally, PromptGen reveals that the CLIP model shows “reporting bias” when used as control, and PromptGen can further de-bias this controlled distribution in an iterative manner.¹

1 Introduction

Generative models learn the underlying high-dimensional data distribution and have achieved promising performance on image synthesis [4, 73, 35, 21]. Though being well praised, they still face two main criticisms. First, since generative models are typically trained in an unsupervised way, they lack controllability, meaning that it is unclear how to sample from a specific region of the space. Second, generative models are prone to inherit the imbalance and bias of training data [58, 30]. For instance, StyleGAN2 is more likely to produce images of white individuals, see Figure 1(e). Previous works have studied these challenges separately, and typical methods include editing of “style” codes [22, 66, 30] and explicit conditions [41, 67]. However, these methods are either model-dependent (i.e., requiring a well-structured style space) or label-intensive (i.e., requiring all training samples to be labeled for explicit conditions), limiting their generality and practical use.

To address the above challenges, this paper advocates a unified formulation, *distributional control of generative models*, which enables controllability by incorporating the knowledge of off-the-shelf models (e.g., CLIP [56], classifiers, or inverse graphics models [16]). Based on this unified view, we propose to learn distributions in the latent space of a pre-trained generative model while keeping the pre-trained weights fixed. Given its conceptual similarity to prompt learning [40, 83, 84, 29], we term our framework as Generative Visual Prompt (PromptGen). PromptGen requires **no training data**, and the only supervision comes from off-the-shelf models that help define the control. Specifically, PromptGen allows the user to define controls using an energy-based model (EBM) approximated

¹Our code is available at <https://github.com/ChenWu98/Generative-Visual-Prompt>.

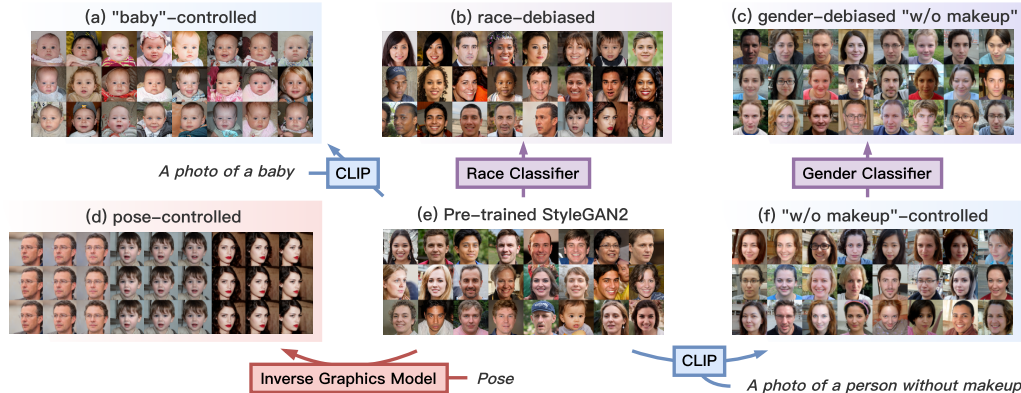


Figure 1: PromptGen uses different off-the-shelf models (e.g., CLIP [56], inverse graphics models, and classifiers) to control the output distribution of a pre-trained and fixed generative model (e.g., StyleGAN2). Colored boxes (e.g., the blue CLIP) indicate the control. See the text for details.

by an invertible neural network (INN). Unlike methods that require optimization at inference in EBM sampling [47, 14, 43, 50], PromptGen samples images in a **feed-forward** manner, which is highly efficient. Moreover, PromptGen **stands alone** at inference, meaning that we can discard the off-the-shelf models, which define the control, after training.

We illustrate the benefits of PromptGen with several experiments. Figure 1 demonstrates our main findings with StyleGAN2 trained on FFHQ [34] without any labels, while we also show results for StyleNeRF [21], diffusion autoencoder [55], and NVAE [73] in Section 4. Figure 1(a) illustrates how we can sample from StyleGAN2 based on text descriptions such as “a photo of a baby”. Figure 1(b) shows that PromptGen can leverage a race classifier (potentially trained on a different dataset) to sample uniformly across all races, de-biasing the pre-trained StyleGAN2. Figure 1(d) illustrates that PromptGen can generate images of the same identity in different poses, guided by a pose regressor.

Finally, it is worth pointing out that PromptGen not only offers **generality** for algorithmic design and **modularity** for control composition, but also enables **iterative** controls when some controls are contingent on others. For instance, one may train a text-controlled distribution and then de-bias this distribution. To achieve this, we view the composition of the INN and the generative model as a new “generative model” to be controlled. Figure 1(f) illustrates that PromptGen reveals “reporting bias” of the CLIP model [56], where “without makeup” is – perhaps surprisingly – positively correlated with female, and Figure 1(c) shows that PromptGen can further mitigate this bias with iterative control.

Table 1: Comparison between methods.

	StyleFlow [1]	PPGM [47] / LACE [50]	Guided DDPM [10]	PromptGen
Arbitrary control (e.g., CLIP)	✗	✓	✓	✓
Low-dimensional latent space	✓	✓	✗	✓
Stands alone at inference	✓	✗	✗	✓
Feed-forward (i.e., no inference-time optim.)	✓	✗	✗	✓
Iterative distributional control	✗	✗	✗	✓

2 Related Work

Over the past few years, generative models have gained the ability to generate images with high visual quality. A few of the most widely used methods include generative adversarial networks (GANs) [18], VAE [38], invertible neural networks [11, 12], and diffusion models [24, 69]. In particular, generative models trained on large amounts of unlabeled data, e.g., BigGAN [4], StyleGANs [34, 35, 33, 62], Glow [37], and diffusion models [49, 10], achieve promising image synthesis results.

Despite their success, controllability and de-biasing are still two fundamental challenges generative models face. For controllability, existing methods include explicit conditioning at training [41, 67]

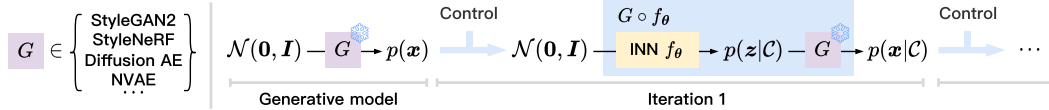


Figure 2: Overview of Generative Visual Prompt (PromptGen). Given a pre-trained generative model G and a control \mathcal{C} , we learn a distribution $p(\mathbf{z}|\mathcal{C})$ in G 's latent space, while keeping G fixed. Each control \mathcal{C} can have multiple components, e.g., $\mathcal{C} = \{\text{text} = \text{"a photo of a baby"}, \text{gender} = \text{male}\}$. PromptGen views the composition $G \circ f_\theta$ as a new “generative model” for iterative control.

or local editing of the learned representation, e.g. “style” codes [1, 54, 44, 66, 78]. For de-biasing, existing methods include local editing of “style” codes [58, 30] and importance sampling for either training [20] or inference [28]. Existing works study these problems separately, each requiring a specific design for the task studied. On the contrary, PromptGen is a unified framework for arbitrary controls defined by off-the-shelf models. The benefits of unifying tasks have been shown by a recent trend of works in multiple research areas across vision and language [57, 46, 61, 79, 59]. Moreover, one can fine-tune generative models [17] for domain adaptation, which is orthogonal to PromptGen: since PromptGen maps a generative model to another generative model (Algorithm 1), fine-tuning can be applied before or after PromptGen training. We leave this exploration to future studies.

Previous methods usually sample from EBM [39] with Markov Chain Monte Carlo (MCMC) [71, 74, 14, 19, 13, 43]. Among them, plug-and-play generative models (PPGMs) [47] and LACE [50] define latent-space EBMs. However, MCMC requires inference-time optimization, which is inefficient and requires the off-the-shelf models to be available at inference; this is also the criticism for diffusion models, regardless of being gradient-guided [49, 65] or not [24, 25, 48]. In contrast, PromptGen achieves efficient, feed-forward sampling. Table 1 shows a comparison with previous methods. Our INN training is similar to that proposed by [52] to sample from physical systems, which is later used by [76] to solve inverse problems with two composed INNs. Notably, [52] and [76] do not leverage a low-dimensional latent space, and [76] requires training a separate model for each image sample. In this paper, we use INN to model arbitrary EBMs for various generative models.

3 Method

Figure 2 illustrates our PromptGen framework. To begin, the user selects a pre-trained generative model G . PromptGen then lets the user specify a control \mathcal{C} as an energy-based model (EBM) $p(\mathbf{z}|\mathcal{C})$. We train an INN f_θ to approximate $p(\mathbf{z}|\mathcal{C})$. PromptGen views the functional composition $G \circ f_\theta$ as a new generative model and can perform iterative control. Algorithm 1 describes the overall procedure.

Algorithm 1: Generative Visual Prompt (PromptGen)

Input: Generative model $G : \mathbb{R}^d \rightarrow \mathcal{X}$

repeat

1. **Input:** control \mathcal{C} of the current iteration
2. Define an EBM $p(\mathbf{z}|\mathcal{C})$ for \mathcal{C} (Section 3.1)
3. Train an INN $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ to approximate $p(\mathbf{z}|\mathcal{C})$ (Section 3.2)
4. $G \leftarrow G \circ f_\theta$

until user stops the iteration

return $G : \mathbb{R}^d \rightarrow \mathcal{X}$

3.1 Latent-Space EBM for Distributional Control

The plug-and-play generative model [47] was first proposed to use a latent-space EBM for controllable image synthesis. If a fixed generative model is used, then theoretically, any image-space EBM can be viewed as a latent-space EBM, as shown by [19] and [50]. We define the latent-space EBM following similar formulation as [47, 19, 50], with some new energy functions. We then extend the formulation to incorporate the moment constraint [7], which was adopted for language modeling [36], but unlike [36], we define the moment constraint in the latent space to accommodate generative vision models.

We define a control \mathcal{C} as M independent properties $\{\mathbf{y}_1, \dots, \mathbf{y}_M\}$, e.g., \mathbf{y}_1 can be a text description and \mathbf{y}_2 can be an attribute. The controllability can be defined as the EBM (detailed in Appendix B.1)

$$p(\mathbf{x}|\mathcal{C}) = \frac{p_{\mathbf{x}}(\mathbf{x})e^{-E_{\mathcal{C}}(\mathbf{x})}}{Z_X}, \quad E_{\mathcal{C}}(\mathbf{x}) = \sum_{i=1}^M \lambda_i E_i(\mathbf{x}, \mathbf{y}_i), \quad Z_X = \int_{\mathbf{x}'} p_{\mathbf{x}}(\mathbf{x}')e^{-E_{\mathcal{C}}(\mathbf{x}')} d\mathbf{x}', \quad (1)$$

which reweights the image prior $p_{\mathbf{x}}(\mathbf{x})$ with energy $E_{\mathcal{C}}(\mathbf{x})$, where images with smaller energy are preferred. Using a pre-trained generative model $G: \mathbb{R}^d \rightarrow \mathcal{X}$ that maps a latent code \mathbf{z} to an image \mathbf{x} , the image prior $p_{\mathbf{x}}(\mathbf{x})$ is defined by sampling a latent code \mathbf{z} from $p_{\mathbf{z}} = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and mapping it to $\mathbf{x} = G(\mathbf{z})$. Appendix B.2 shows that this EBM is equivalent to the latent-space EBM

$$p(\mathbf{z}|\mathcal{C}) = \frac{p_{\mathbf{z}}(\mathbf{z})e^{-E_{\mathcal{C}}(G(\mathbf{z}))}}{Z}, \quad E_{\mathcal{C}}(\mathbf{x}) = \sum_{i=1}^M \lambda_i E_i(\mathbf{x}, \mathbf{y}_i), \quad Z = \int_{\mathbf{z}'} p_{\mathbf{z}}(\mathbf{z}')e^{-E_{\mathcal{C}}(G(\mathbf{z}'))} d\mathbf{z}'. \quad (2)$$

Latent-space EBM allows us to use any off-the-shelf model to specify the control. The following are some examples that are discussed in this paper (explained in Appendix B.1):

Classifier energy: Given a classifier $P(\cdot|\mathbf{x})$ and the target class a that we want to sample images from, we define the classifier energy as $E_{\text{classifier}}(\mathbf{x}, a) = -\log P(a|\mathbf{x})$.

CLIP energy: Using the CLIP model [56], we define the CLIP energy as the cosine distance between the embeddings of the image and the text \mathbf{t} , averaged over L differentiable augmentations [82, 44]:

$$E_{\text{CLIP}}(\mathbf{x}, \mathbf{t}) = \frac{1}{L} \sum_{l=1}^L \left(1 - \cos \left\langle \text{CLIP}_{\text{img}}(\text{DiffAug}_l(\mathbf{x})), \text{CLIP}_{\text{text}}(\mathbf{t}) \right\rangle \right). \quad (3)$$

Inverse graphics energy: Given an inverse graphics model, $f_{\mathcal{X} \rightarrow \mathcal{P}}$, which infers image parameters (e.g., pose and expression), and the target parameters $\boldsymbol{\rho}$, we define the inverse graphics energy as

$$E_{\text{inv-graphics}}(\mathbf{x}, \boldsymbol{\rho}) = d \left\langle f_{\mathcal{X} \rightarrow \mathcal{P}}(\mathbf{x}), \boldsymbol{\rho} \right\rangle^2, \quad (4)$$

where $d\langle \cdot, \cdot \rangle$ is the geodesic distance between the inferred parameters and the target parameters.

Moment constraint: Some controls cannot be *directly* defined by off-the-shelf models, and the moment constraint [7, 36] is one of them. Given a mapping $\gamma: \mathcal{X} \rightarrow \mathbb{R}^K$ (e.g., γ can be a classifier that outputs the probability simplex), the moment constraint defines the target distribution $p(\mathbf{x}|\mathcal{C})$ as

$$p(\mathbf{x}|\mathcal{C}) = \underbrace{\arg \min_{p(\mathbf{x}|\mathcal{C})} \mathbb{D}_{\text{KL}}(p(\mathbf{x}|\mathcal{C}) \| p_{\mathbf{x}}(\mathbf{x}))}_{\text{Deviation from the pre-trained distribution}}, \quad \text{s.t.} \quad \underbrace{\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathcal{C})} [\gamma(\mathbf{x})]}_{\text{Moment constraint}} = \boldsymbol{\mu}, \quad (5)$$

where $\boldsymbol{\mu}$ is the user-specified vector. For example, if we want to generate images that are uniformly distributed across races, we may use a race classifier as γ and define $\boldsymbol{\mu} = (|\mathcal{A}|^{-1}, \dots, |\mathcal{A}|^{-1})$ where \mathcal{A} is the set of races. In this paper, we generalize the moment constraint to the latent space, and approximate the above objective as (detailed in Appendix B.6):

$$p(\mathbf{z}|\mathcal{C}) = \frac{p_{\mathbf{z}}(\mathbf{z}) \exp \left(\hat{\boldsymbol{\beta}}^\top \gamma(G(\mathbf{z})) \right)}{Z}, \quad Z = \int_{\mathbf{z}'} p_{\mathbf{z}}(\mathbf{z}') \exp \left(\hat{\boldsymbol{\beta}}^\top \gamma(G(\mathbf{z}')) \right) d\mathbf{z}', \quad (6)$$

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(N)} \stackrel{\text{i.i.d.}}{\sim} p_{\mathbf{z}}(\mathbf{z}), \mathbf{x}^{(j)} = G(\mathbf{z}^{(j)})} \left\| \frac{\sum_{j=1}^N \exp(\boldsymbol{\beta}^\top \gamma(\mathbf{x}^{(j)})) \gamma(\mathbf{x}^{(j)})}{\sum_{j'=1}^N \exp(\boldsymbol{\beta}^\top \gamma(\mathbf{x}^{(j')}))} - \boldsymbol{\mu} \right\|_2^2. \quad (7)$$

3.2 Approximating EBM with Invertible Neural Network

Given the functional form of the EBM $p(\mathbf{z}|\mathcal{C})$, our next step is to approximate it with an efficient sampling network. To achieve this, we train a distribution $p_{\boldsymbol{\theta}}(\mathbf{z})$ that minimizes the KL divergence $\mathbb{D}_{\text{KL}}(p_{\boldsymbol{\theta}}(\mathbf{z}) \| p(\mathbf{z}|\mathcal{C}))$. Estimating $\mathbb{D}_{\text{KL}}(p_{\boldsymbol{\theta}}(\mathbf{z}) \| p(\mathbf{z}|\mathcal{C}))$ requires easily sampling $\mathbf{z} \sim p_{\boldsymbol{\theta}}$ and tractably computing $p_{\boldsymbol{\theta}}(\mathbf{z})$. Inspired by [52], we model $p_{\boldsymbol{\theta}}$ with an INN $f_{\boldsymbol{\theta}}$ that defines a bijection $\mathbf{z} = f_{\boldsymbol{\theta}}(\boldsymbol{\epsilon})$, which has two merits besides the invertibility: (1) one can easily sample \mathbf{z} by sampling $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and mapping it to $\mathbf{z} = f_{\boldsymbol{\theta}}(\boldsymbol{\epsilon})$, and (2) $p_{\boldsymbol{\theta}}(\mathbf{z})$ has a closed-form solution:

$$\log p_{\boldsymbol{\theta}}(\mathbf{z}) = \log \mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \mathbf{I}) - \log \left| \det \left(\frac{\partial f_{\boldsymbol{\theta}}}{\partial \boldsymbol{\epsilon}} \right) \right|, \quad \mathbf{z} = f_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}). \quad (8)$$

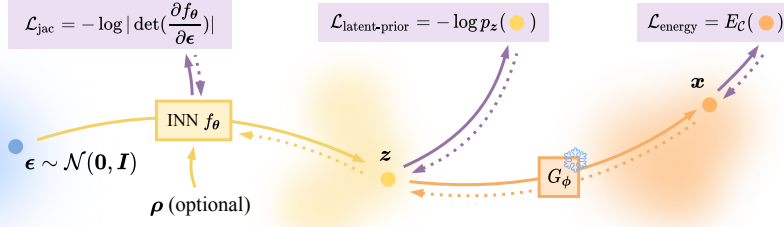


Figure 3: Illustration for Algorithm 2. In the forward pass (solid curves), we sample $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, map ϵ to latent code $z = f_\theta(\epsilon)$ with an INN f_θ , and map z to an image $x = G(z)$ with a fixed generative model G . Dashed curves show the gradients. Details are provided in Section 3.2.

Algorithm 2: Approximating Latent-Space EBM with INN

while not converged do

1. Sample $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 2. Map ϵ to latent code $z = f_\theta(\epsilon)$
 3. Map z to an image $x = G(z)$
 4. Optimize θ with gradient $\nabla_\theta \left(-\log \left| \det \left(\frac{\partial f_\theta}{\partial \epsilon} \right) \right| - \log p_z(z) + E_C(x) \right)$
-

Based on these properties of INN, we can rewrite our KL divergence objective $\mathbb{D}_{\text{KL}}(p_\theta(z) \| p(z|\mathcal{C}))$ as (full derivations in Appendix B.3) the following form:

$$\begin{aligned} \mathbb{D}_{\text{KL}}(p_\theta(z) \| p(z|\mathcal{C})) &= \mathbb{E}_{z \sim p_\theta(z), x = G(z)} \left[\log \frac{p_\theta(z)}{p_z(z) e^{-E_C(x)} / Z} \right] \\ &= \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), z = f_\theta(\epsilon), x = G(z)} \left[-\log \left| \det \left(\frac{\partial f_\theta}{\partial \epsilon} \right) \right| - \log p_z(z) + E_C(x) \right] - \mathbb{H}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} + \log Z. \end{aligned} \quad (9)$$

Since $\mathbb{H}_{\mathcal{N}(\mathbf{0}, \mathbf{I})}$ and $\log Z$ are independent of θ , our training objective becomes

$$\arg \min_{\theta} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), z = f_\theta(\epsilon), x = G(z)} \left[\underbrace{-\log \left| \det \left(\frac{\partial f_\theta}{\partial \epsilon} \right) \right|}_{\mathcal{L}_{\text{jac}}} - \underbrace{\log p_z(z)}_{\mathcal{L}_{\text{latent-prior}}} + \underbrace{E_C(x)}_{\mathcal{L}_{\text{energy}}} \right]. \quad (10)$$

Figure 3 gives an illustration of our process, and Algorithm 2 describes the algorithmic details.

PromptGen in a class-embedding space (Figure 5(e)) Previous works [4, 62] have shown that *class conditioning* boosts generative models’ performances on ImageNet [60]. Specifically, class-conditioned generative models map a latent code z and a class embedding y to $x = G(z, y)$. To extend PromptGen to these models, we train an INN h_θ to map $\xi \sim \mathcal{N}(\mu, \sigma^2 \mathbf{I})$ to $y = h_\theta(\xi)$, where μ and σ are the mean and standard deviation of G ’s class embeddings. The motivation for defining the distribution of ξ as $\mathcal{N}(\mu, \sigma^2 \mathbf{I})$ but not $\mathcal{N}(\mathbf{0}, \mathbf{I})$ is that we want the learned INN h_θ to be volume-preserving, which is easier to train. The training objective is to minimize $\mathbb{D}_{\text{KL}}(p_\theta(z, y) \| p(z, y|\mathcal{C}))$, which is equivalent to (full derivations in Appendix B.4)

$$\begin{aligned} \arg \min_{\theta} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \xi \sim \mathcal{N}(\mu, \sigma^2 \mathbf{I}), z = f_\theta(\epsilon), y = h_\theta(\xi), x = G(z, y)} \left[-\log \left| \det \left(\frac{\partial f_\theta}{\partial \epsilon} \right) \right| - \log p_z(z) \right. \\ \left. - \log \left| \det \left(\frac{\partial h_\theta}{\partial \xi} \right) \right| - \log p_y(y) + E_C(x) \right], \end{aligned} \quad (11)$$

where $p_y(y) = \mathcal{N}(y|\mu, \sigma^2 \mathbf{I})$ is the estimated class-embedding distribution.

PromptGen with conditional INN To generalize the control to continuous values, e.g., the scene parameters ρ in Eq. (4), we condition the INN on ρ . The condition is modeled by replacing $z = f_\theta(\epsilon)$ with $z = f_\theta(\epsilon, \rho)$. Space limited, we provide details and derivations in Appendix B.5. This extension results in a similar architecture to StyleFlow [1], but StyleFlow [1] uses MLE training on image-label pairs and is only applicable to explicit condition, which is not capable of modeling EBMs. When the condition is modeled by the equivariant operations proposed by [77], PromptGen also satisfies the homomorphism property [77] in the latent space.



Figure 4: Image synthesis based on text description, guided by the CLIP model. As the pre-trained generative model, we use StyleGAN2 trained on FFHQ 1024² [34] with truncation $\psi = 0.7$. The text description used in this experiment is a photo of a baby. StyleCLIP requires optimization at inference, while StyleGAN2-NADA and PromptGen do not. All images are resized for visualization.

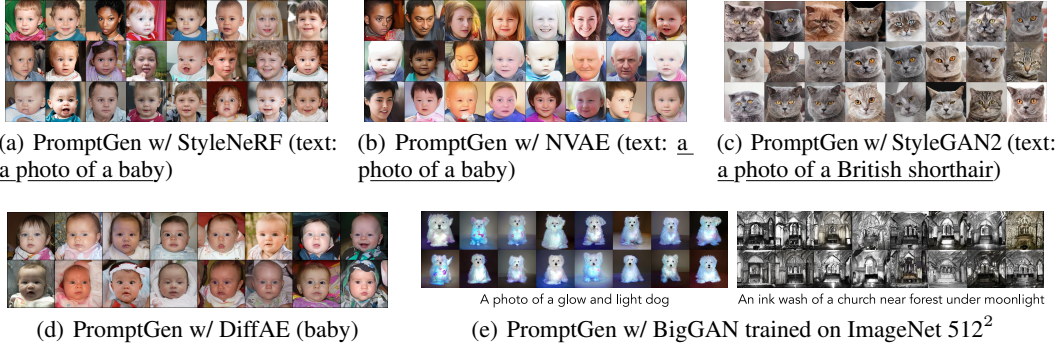


Figure 5: PromptGen is applicable to different generative models in various domains. Figure 5(a), Figure 5(d), and Figure 5(b) are PromptGen applied to StyleNeRF [21], diffusion autoencoder [55], and NVAE [73]. Figure 5(c) is PromptGen applied to StyleGAN2 [35] trained on AFHQ-Cats [6]. Figure 5(e) is the extension to the embedding space of BigGAN [4] on ImageNet [60]. All images are resized for visualization. See Appendix C for results on more datasets and text descriptions.

4 Experiments

This section describes the experimental validation of PromptGen. See Appendix A for experiments on synthetic data and Appendix C, D, E, and F for additional experiments on images and 3D meshes. Additional experimental details are provided in Appendix B.8.

4.1 Image Synthesis based on Text Description

This experiment illustrates the capability of PromptGen to sample images from a generative model driven by a text description t using the pre-trained CLIP model [56] (the ViT-B/32 version). We used the CLIP energy from Eq. (3), with text descriptions such as a photo of a baby.

Figure 4 shows a comparison between PromptGen and two previous CLIP-guided image generation methods, StyleCLIP [54] and StyleGAN2-NADA [17]. We observe that PromptGen generates diverse and high-quality images of babies, while StyleCLIP struggles in controllability and image quality, and StyleGAN2-NADA generates baby-like adults.² These results show that (1) locally editing the latent code (i.e., StyleCLIP) is not always an effective method for controlling generative models (e.g., not all images’ latent code can be locally edited into a baby); (2) domain adaptation (i.e., StyleGAN2-NADA) is not effective in seeking modes in a generative model. Figure 5(a), Figure 5(b), and Figure 5(d) show how PromptGen applies to StyleNeRF [21], NVAE [73], and diffusion autoencoder (DiffAE; [55]). Although these generative models are all trained on FFHQ, baby images sampled from them have distinct characteristics. Figure 5(c) shows the ability of PromptGen to generate cats of a particular species. Figure 5(e) shows the extension (Section 3.2) of our PromptGen to the class-embedding space of BigGAN [4], a class-conditional GAN trained on ImageNet [60]. We observe that PromptGen helps BigGAN generate images with complex text descriptions, which are out of BigGAN’s training image distribution; however, the diversity seems to be limited in this extension.

²Since we care about distributions, *none of the images in this paper are cherry- or lemon-picked.*

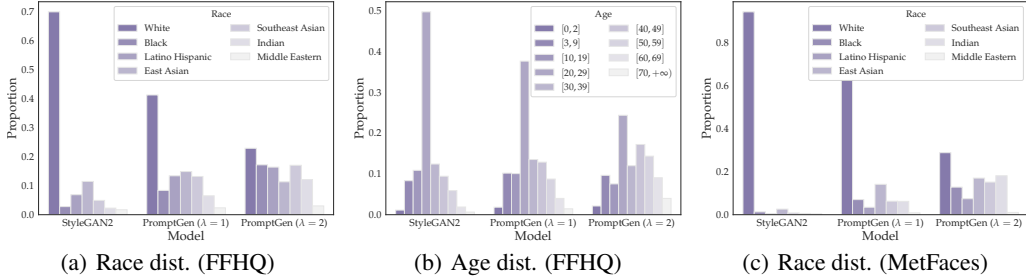


Figure 6: With the moment constraint, PromptGen de-biases StyleGAN2 (on FFHQ and MetFaces 1024^2 , truncation $\psi = 0.7$). See Appendix D for image samples and Table 5 for quantitative results.

Table 2: Comparison with baselines for de-biasing binary attributes and their correlations. Following [30], we use classifiers on CelebA [45]. Baseline performances are copied from [30]. PromptGen has competitive performance, even in the cases where FairStyle achieves nearly perfect performance.

	FFHQ (binary attributes)					
	$\mathbb{D}_{\text{KL}}^{\text{gender}} \downarrow$	$\mathbb{D}_{\text{KL}}^{\text{eyeglasses}} \downarrow$	$\mathbb{D}_{\text{KL}}^{\text{blond hair}} \downarrow$	$\mathbb{D}_{\text{KL}}^{\text{age+gender}} \downarrow$	$\mathbb{D}_{\text{KL}}^{\text{age+eyeglasses}} \downarrow$	$\mathbb{D}_{\text{KL}}^{\text{gender+eyeglasses}} \downarrow$
FFHQ (real data)	0.015	0.186	–	0.246	0.355	0.242
StyleGAN2 [35]	0.018	0.180	–	0.279	0.384	0.250
StyleFlow [1]	0.023	0.061	–	0.214	0.162	0.121
FairGen [70]	4.21×10^{-4}	7.07×10^{-4}	–	0.0373	0.0330	0.00185
FairStyle [30]	3.20×10^{-7}	0	–	0.0257	0.0157	0.000241
PromptGen (ours)	1.71×10^{-5}	1.72×10^{-5}	0.0008	0.000558	0.000415	0.000628

4.2 De-Biasing Pre-Trained Generative Models

An important problem in generative models is to generate fair distributions w.r.t a set of attributes of interest. For instance, Figure 6(a) shows that StyleGAN2 generates images with bias across races and ages. PromptGen de-biases StyleGAN2 models trained on FFHQ 1024^2 and MetFaces 1024^2 [32] in terms of *categorical* attributes, using the moment constraint defined in Eq. (5) and Eq. (6). As control, we used a classifier trained on FairFace 224^2 [31] as γ . We defined $\mu = (|\mathcal{A}|^{-1}, \dots, |\mathcal{A}|^{-1})$, where \mathcal{A} is the set of races. Similar to the energy weights λ_i defined in Eq. (2), we propose to rescale the trained $\hat{\beta}$ as $\lambda\hat{\beta}$. Figure 6 shows that PromptGen de-biases the race and age effectively.

Existing de-biasing baselines consider *binary* attributes [30]. For a fair comparison with them, we adopted their setting to use binary classifiers trained on CelebA [45] for de-biasing and evaluation. Since classifiers trained on CelebA also suffer from the spurious correlation between attributes, we did *not* use the moment constraint for this experiment. Instead, since PromptGen allows conditional image generation with the classifier energy, we generated the same number of samples conditioned on each attribute or attribute combination. Table 2 shows that PromptGen has competitive performance on de-biasing for attributes and attribute combinations.

4.3 Pose-Guided Face Synthesis

With an inverse graphics model, PromptGen can control the pose of faces generated by StyleGAN2. We used the DECA model [16], which infers the parameters of FLAME [42], a parametric facial graphics model. We set $\rho \in \text{SO}(3)$ as FLAME’s three neck poses and used the conditional INN extension introduced in Section 3.2. To enable generating different poses of the same identity (ID), we propose ID energy using the IR-SE50 model [8]. Specifically, given a canonical pose ρ_0 , we define $z_0 = f_\theta(\epsilon, \rho_0)$, and the ID energy is defined as (detailed in Appendix B.7)

$$E_{\text{ID}}(x_0, x) = 1 - \cos \langle R(x_0), R(x) \rangle, \quad x_0 = G(z_0), x = G(z), \quad (12)$$

where R is the IR-SE50 model [8] that computes face embeddings. Figure 7 shows that PromptGen generates faces of the same ID in different poses, even without being explicitly trained with poses as conditions. We computed the FID score [23] of each model using Clean-FID [53]; following [35],



Figure 7: Using inverse graphics model DECA [16] and scene parameters $\rho_{\text{pose}} \in \text{SO}(3)$, PromptGen controls the pose of StyleGAN2 while preserving the identity. All images are resized for visualization.

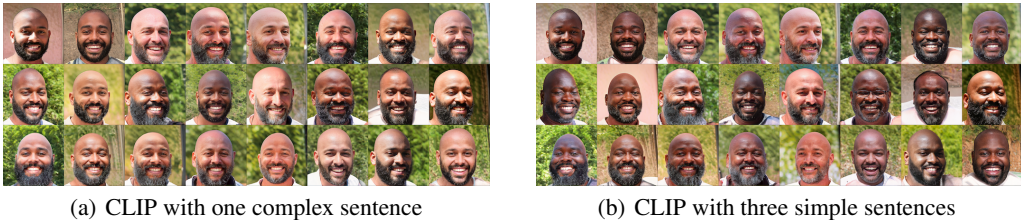
Table 3: Pose-controlled face generation. In this experiment, PromptGen uses StyleGAN2 as the pre-trained generative model. Results for GRAF [64], pi-GAN [5], GIRAFFE [51], and StyleNeRF [21] are from [21]; results for DiscoFaceGAN (DFG) [9] and GAN-Control [67] are from [67].

	StyleGAN2	GRAF	pi-GAN	GIRAFFE	StyleNeRF	DFG	GAN-Control	PromptGen
Resolution	1024 ²	256 ²	256 ²	256 ²	1024 ²	256 ²	512 ²	1024 ²
Pose	✗	✓	✓	✓	✓	✓	✓	✓
FID↓	3	71	85	35	8	13	6	4
Dist. w/ same ID↓	–	–	–	–	–	0.83	0.68	0.45
Dist. w/ diff. ID	–	–	–	–	–	1.73	1.90	1.37

we did not use the truncation trick when computing the FID score. Table 3 shows that PromptGen outperforms existing models in terms of the FID score. Following [67], we then reported the average IR-SE50 [8] embedding distances for images with the same ID and with different IDs. Results show that PromptGen achieves the best ID preservation, with a slight sacrifice of ID diversity.

4.4 Decomposing Complex Control via Energy Composition

Energy composition $E_C(\mathbf{x}) = \sum_{i=1}^M \lambda_i E_i(\mathbf{x}, \mathbf{y}_i)$ allows us to decompose controls into simple ones. We show that $\mathbf{y}_1 = \text{photo of a bald black man with beard}$ ($\lambda_1 = 6000$) is not successful (Figure 8(a)); by decomposing it as $\mathbf{y}_1 = \text{photo of a bald man}$ ($\lambda_1 = 1500$), $\mathbf{y}_2 = \text{photo of a black man}$ ($\lambda_2 = 3000$), and $\mathbf{y}_3 = \text{photo of a man with beard}$ ($\lambda_3 = 1500$), we have better control (Figure 8(b)).



(a) CLIP with one complex sentence

(b) CLIP with three simple sentences

Figure 8: Decomposing complex controls (e.g., photo of a bald black man with beard) into simpler ones improves control performance. See details in Section 4.4.

4.5 Iterative Distributional Control via Functional Composition

This section discusses an interesting bias that PromptGen reveals about the CLIP model. Figure 9(a) shows images generated by PromptGen (with StyleGAN2) with the CLIP model and the text description a photo of a person without makeup, where more females are generated than males. This bias should not be attributed to *image pre-training data* since images contain a bias in the opposite direction, i.e., men are less likely to have makeup. We argue that this bias should be explained by CLIP having learned a “reporting bias” in *vision-language pre-training data*: people are more likely to say “a person without makeup” when the person is a female (detailed analysis in Appendix F).

Besides revealing the above “reporting bias”, PromptGen can also mitigate this bias via an iterative control, enabled by the functional composition in Algorithm 1. Specifically, in the second iteration, we de-biased the gender distribution of $G \circ f_\theta$ instead of G , where f_θ is the INN learned for the text control. For de-biasing, we used the moment constraint with $\hat{\beta}$ trained for $G \circ f_\theta$. Figure 9(b) shows that females and males are uniformly distributed after the moment constraint in the second iteration.



(a) PromptGen (iteration 1). Image synthesis with text description a photo of a person without makeup. Gender distribution: female: 81.6%; male: 18.4%.



(b) PromptGen (iteration 2). Gender de-biasing ($\lambda = 2$) for the distribution learned in iteration 1. Gender distribution: female: 49.3%; male: 50.7%.

Figure 9: A curious case of the CLIP model: with description a photo of a person without makeup, PromptGen generates more female images than male, showing that CLIP learns a “reporting bias”. Besides revealing this “reporting bias”, PromptGen can also mitigate this bias via an iterative control, allowing us to de-bias the text-controlled distribution. All images are resized for visualization.

4.6 Inference Latency

We compared our PromptGen and the plug-and-play generative model (PPGM) [47] in terms of the inference latency. PPGM uses Langevin dynamics [74] to optimize over the latent space at inference, while PromptGen samples images in a feed-forward manner. In this experiment, we used PPGM and our PromptGen to approximate the same EBM with CLIP energy. Inference times were estimated on an NVIDIA RTX A4000 GPU. Table 4 shows that our PromptGen has the highest performance and efficiency. Compared to our PromptGen, PPGM requires $100\times$ inference time to achieve comparable results. Moreover, learning f_θ in PromptGen requires training-time optimization, and this amortized optimization is useful when one wants to reuse a controlled distribution many times.

Table 4: Comparison between PromptGen and PPGM, when approximating the *same* EBM with CLIP energy. n is the number of inference-time optimization steps used by PPGM. All models used StyleGAN2 as the pre-trained generative model. The text descriptions are a photo of a {baby, boy, girl}.

	PPGM ($n = 10$)	PPGM ($n = 50$)	PromptGen (ours)
CLIP energy (baby)↓	0.7327	0.7134	0.7038
CLIP energy (girl)↓	0.7257	0.7184	0.7199
CLIP energy (boy)↓	0.7263	0.7114	0.7081
Inference time per sample (sec.)↓	4.4	21.5	0.2
Back-propagation through CLIP at inference	10	50	0
When is it equal to the EBM?	$n \rightarrow \infty$	$n \rightarrow \infty$	$\mathbb{D}_{\text{KL}}(p_\theta(z) p(z C)) = 0$

5 Conclusions and Future Work

This paper proposes PromptGen, a unified framework to learn latent distributions for distributional control of pre-trained generative models. PromptGen leverages the knowledge of various off-the-shelf models and, unlike previous methods, it does not require the availability of these models at inference. PromptGen can sample images in a feed-forward manner, which is more efficient than methods that require optimization at inference. PromptGen offers the generality for algorithmic design and modularity for control composition, and it also enables iterative controls. Experiments validate that PromptGen applies to various generative models (StyleGAN2, StyleNeRF, NVAE), control types (continuous, discrete, and moment constraint), off-the-shelf models (CLIP, classifiers, and inverse graphics models), and data domains (faces, churches, animals, ImageNet, and landscapes).

Limitation and future work: We provide an error analysis in Appendix F and a discussion on societal impact in Appendix G. PromptGen is restricted by the pre-trained generative model’s coverage [3, 27]. PromptGen focuses on mode-seeking and -reweighting instead of domain adaptation, and it is possible to combine PromptGen and domain adaptation of generative models [17]. Also, PromptGen depends on the off-the-shelf models that provide knowledge about the control (Appendix F). In subsequent work, it can be beneficial to explore how to learn energy functions [15, 26] (besides our moment constraint used for de-biasing), which may provide fine-grained control with less bias.

Acknowledgments and Disclosure of Funding

The authors would like to thank the anonymous reviewers, Zoltán Ádám Milacski, Jianchun Chen, and Shubhra Aich for their valuable feedback on drafts of this paper.

References

- [1] Rameen Abdal, Peihao Zhu, Niloy Jyoti Mitra, and Peter Wonka. StyleFlow: Attribute-conditioned exploration of StyleGAN-generated images using conditional continuous normalizing flows. *TOG*, 2021.
- [2] Lynton Ardizzone, Carsten Lüth, Jakob Kruse, Carsten Rother, and U. Köthe. Guided image generation with conditional invertible neural networks. *ArXiv*, 2019.
- [3] David Bau, Jun-Yan Zhu, Jonas Wulff, William S. Peebles, Hendrik Strobelt, Bolei Zhou, and Antonio Torralba. Seeing what a GAN cannot generate. *ICCV*, 2019.
- [4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *ICLR*, 2019.
- [5] Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-GAN: Periodic implicit generative adversarial networks for 3D-aware image synthesis. *CVPR*, 2021.
- [6] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. StarGAN v2: Diverse image synthesis for multiple domains. *CVPR*, 2020.
- [7] Imre Csiszár and Paul C. Shields. Information theory and statistics: A tutorial. *Found. Trends Commun. Inf. Theory*, 2004.
- [8] Jiankang Deng, J. Guo, and Stefanos Zafeiriou. ArcFace: Additive angular margin loss for deep face recognition. *CVPR*, 2019.
- [9] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and controllable face image generation via 3D imitative-contrastive learning. *CVPR*, 2020.
- [10] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. *NeurIPS*, 2021.
- [11] Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: Non-linear independent components estimation. *ICLR, Workshop Track Proceedings*, 2015.
- [12] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. *ICLR*, 2017.
- [13] Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models. *NeurIPS*, 2020.
- [14] Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *NeurIPS*, 2019.
- [15] Jesse H. Engel, Matthew D. Hoffman, and Adam Roberts. Latent constraints: Learning to generate conditionally from unconditional generative models. *ICLR*, 2018.
- [16] Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. Learning an animatable detailed 3D face model from in-the-wild images. *TOG*, 2021.
- [17] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. StyleGAN-NADA: CLIP-guided domain adaptation of image generators. *SIGGRAPH*, 2022.
- [18] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. *NIPS*, 2014.
- [19] Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. *ICLR*, 2020.

- [20] Aditya Grover, Kristy Choi, Rui Shu, and Stefano Ermon. Fair generative modeling via weak supervision. *ICML*, 2020.
- [21] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. StyleNeRF: A style-based 3D aware generator for high-resolution image synthesis. *ICLR*, 2022.
- [22] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. GANSpace: Discovering interpretable GAN controls. *NeurIPS*, 2020.
- [23] 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. *NIPS*, 2017.
- [24] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020.
- [25] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *NeurIPS Workshop*, 2021.
- [26] Zhiting Hu, Zichao Yang, Ruslan Salakhutdinov, Xiaodan Liang, Lianhui Qin, Haoye Dong, and Eric P. Xing. Deep generative models with learnable knowledge constraints. *NeurIPS*, 2018.
- [27] Minyoung Huh, Richard Zhang, Jun-Yan Zhu, Sylvain Paris, and Aaron Hertzmann. Transforming and projecting images into class-conditional generative networks. *ECCV*, 2020.
- [28] Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. MaGNET: Uniform sampling from deep generative network manifolds without retraining. *ICLR*, 2022.
- [29] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge J. Belongie, Bharath Hariharan, and Ser Nam Lim. Visual prompt tuning. *ArXiv*, 2022.
- [30] Cemre Karakas, Alara Dirik, Eylül Yalçınkaya, and Pinar Yanardag. FairStyle: Debiasing StyleGAN2 with style channel manipulations. *ArXiv*, 2022.
- [31] Kimmo Kärkkäinen and Jungseock Joo. FairFace: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. *WACV*, 2021.
- [32] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. *NeurIPS*, 2020.
- [33] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *NeurIPS*, 2021.
- [34] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *CVPR*, 2019.
- [35] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. *CVPR*, 2020.
- [36] Muhammad Khalifa, Hady Elsahar, and Marc Dymetman. A distributional approach to controlled text generation. *ICLR*, 2021.
- [37] Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *NeurIPS*, 2018.
- [38] Diederik P. Kingma and Max Welling. Auto-encoding variational Bayes. *ICLR*, 2014.
- [39] Yann LeCun, Sumit Chopra, Raia Hadsell, Aurelio Ranzato, and Fu Jie Huang. A tutorial on energy-based learning. *Tutorial*, 2006.
- [40] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *EMNLP*, 2021.
- [41] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip H. S. Torr. Controllable text-to-image generation. *NeurIPS*, 2019.

- [42] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *TOG*, 2017.
- [43] Nan Liu, Shuang Li, Yilun Du, Joshua B. Tenenbaum, and Antonio Torralba. Learning to compose visual relations. *NeurIPS*, 2021.
- [44] Xingchao Liu, Chengyue Gong, Lemeng Wu, Shujian Zhang, Haoran Su, and Qiang Liu. Fuse-Dream: Training-free text-to-image generation with improved CLIP+GAN space optimization. *ArXiv*, 2021.
- [45] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. *ICCV*, 2015.
- [46] Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 12-in-1: Multi-task vision and language representation learning. *CVPR*, 2020.
- [47] Anh M Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, and Jason Yosinski. Plug & play generative networks: Conditional iterative generation of images in latent space. *CVPR*, 2017.
- [48] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. *ArXiv*, 2021.
- [49] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *ICML*, 2021.
- [50] Weili Nie, Arash Vahdat, and Anima Anandkumar. Controllable and compositional generation with latent-space energy-based models. *NeurIPS*, 2021.
- [51] Michael Niemeyer and Andreas Geiger. GIRAFFE: Representing scenes as compositional generative neural feature fields. *CVPR*, 2021.
- [52] Frank Noé, Jonas Köhler, and Hao Wu. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science*, 2019.
- [53] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in GAN evaluation. *CVPR*, 2022.
- [54] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and D. Lischinski. StyleCLIP: Text-driven manipulation of StyleGAN imagery. *ICCV*, 2021.
- [55] Konpat Preechakul, Nattanat Chatthee, Suttisak Wizadwongsa, and Supasorn Suwajanakorn. Diffusion autoencoders: Toward a meaningful and decodable representation. *CVPR*, 2022.
- [56] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *ICML*, 2021.
- [57] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text Transformer. *JMLR*, 2020.
- [58] Vikram V. Ramaswamy, Sunnie S. Y. Kim, and Olga Russakovsky. Fair attribute classification through latent space de-biasing. *CVPR*, 2021.
- [59] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Manfred Otto Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. A generalist agent. *ArXiv*, 2022.
- [60] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet large scale visual recognition challenge. *IJCV*, 2015.

- [61] Chitwan Saharia, William Chan, Huiwen Chang, Chris A. Lee, Jonathan Ho, Tim Salimans, David J. Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. *ArXiv*, 2021.
- [62] Axel Sauer, Katja Schwarz, and Andreas Geiger. StyleGAN-XL: Scaling StyleGAN to large diverse datasets. *SIGGRAPH*, 2022.
- [63] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. LAION-400M: Open dataset of CLIP-filtered 400 million image-text pairs. *ArXiv*, 2021.
- [64] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. *NeurIPS*, 2020.
- [65] Vikash Sehwal, Caner Hazirbas, Albert Gordo, Firat Ozgenel, and Cristian Canton Ferrer. Generating high fidelity data from low-density regions using diffusion models. *CVPR*, 2022.
- [66] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. InterFaceGAN: Interpreting the disentangled face representation learned by GANs. *TPAMI*, 2022.
- [67] Alon Shoshan, Nadav Bhonker, Igor Kviatkovsky, and Gérard Medioni. GAN-Control: Explicitly controllable GANs. *ICCV*, 2021.
- [68] Ivan Skorokhodov, Grigorii Sotnikov, and Mohamed Elhoseiny. Aligning latent and image spaces to connect the unconnectable. *ICCV*, 2021.
- [69] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *NeurIPS*, 2019.
- [70] Shuhan Tan, Yujun Shen, and Bolei Zhou. Improving the fairness of deep generative models without retraining. *ArXiv*, 2020.
- [71] Tijmen Tieleman. Training restricted Boltzmann machines using approximations to the likelihood gradient. *ICML*, 2008.
- [72] Cristian Vaccari and Andrew Chadwick. Deepfakes and disinformation: Exploring the impact of synthetic political video on deception, uncertainty, and trust in news. *Social Media + Society*, 6, 2020.
- [73] Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. *NeurIPS*, 2020.
- [74] Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient Langevin dynamics. *ICML*, 2011.
- [75] Mika Westerlund. The emergence of Deepfake technology: A review. *Technology Innovation Management Review*, 9(11), 2019.
- [76] Jay Whang, Erik M. Lindgren, and Alexandros G. Dimakis. Composing normalizing flows for inverse problems. In *ICML*, 2021.
- [77] Daniel E. Worrall, Stephan J. Garbin, Daniyar Turmukhambetov, and Gabriel J. Brostow. Interpretable transformations with encoder-decoder networks. *ICCV*, 2017.
- [78] Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. GAN inversion: A survey. *ArXiv*, 2021.
- [79] Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga, Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I. Wang, Victor Zhong, Bailin Wang, Chengzu Li, Connor Boyle, Ansong Ni, Ziyu Yao, Dragomir Radev, Caiming Xiong, Lingpeng Kong, Rui Zhang, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. UnifiedSKG: Unifying and multi-tasking structured knowledge grounding with text-to-text language models. *ArXiv*, 2022.

- [80] Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun Cao. FaceScape: A large-scale high quality 3D face dataset and detailed riggable 3D face prediction. *CVPR*, 2020.
- [81] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop. *ArXiv*, 2015.
- [82] Shengyu Zhao, Zhijian Liu, Ji Lin, Jun-Yan Zhu, and Song Han. Differentiable augmentation for data-efficient GAN training. *NeurIPS*, 2020.
- [83] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *ArXiv*, 2021.
- [84] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. *CVPR*, 2022.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
In Section 5 and Appendix F.
 - (b) Did you describe the limitations of your work? [Yes]
In Section 5 and Appendix F.
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes]
In Section 5.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
In Section 3 and Appendix B.
 - (b) Did you include complete proofs of all theoretical results? [Yes]
In Section 3 and Appendix B.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
We will make our code publicly available on GitHub, with instruction for each experiment. We made great efforts to ensure that most experiments will be run in one line of command.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
In Appendix B.8; other training details will be included in the code release.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
In Appendix B.8.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
We have cited the original papers.
 - (b) Did you mention the license of the assets? [Yes]
License will be included in the code.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
We do not have new assets.
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes]
All data we use are published datasets.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
All data we use are published datasets.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
We do not have user study.
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
We do not have user study.
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]
We do not have user study.

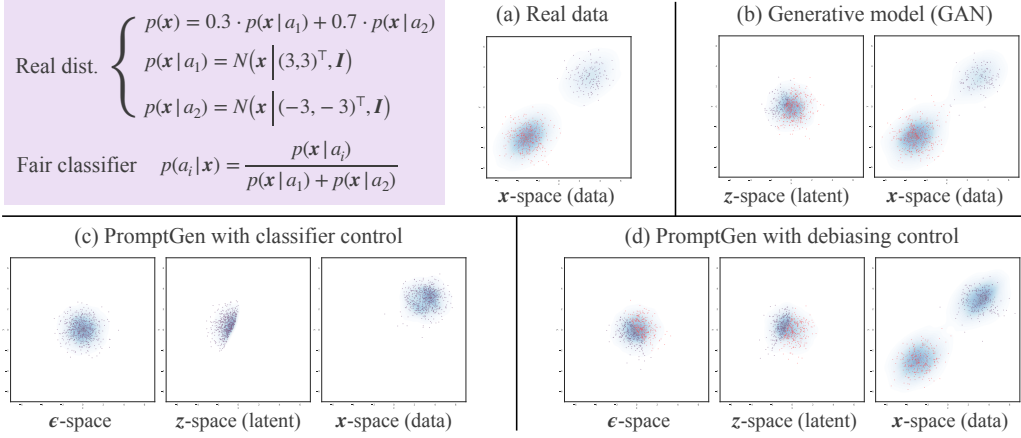


Figure 10: PromptGen on synthetic data. (a) The real distribution is a Gaussian mixture distribution biased towards one mode. We derive a closed-form fair classifier. (b) A GAN trained on this real distribution. (c) With the fair classifier as control, PromptGen learns a distribution concentrated in a specific region. (d) We derive the moment constraint following Section 3.1 and Section 4.2, and PromptGen learns to up-weight the under-represented regions.

A PromptGen on Synthetic Data

We demonstrate the behavior of our PromptGen with two-dimensional synthetic data, using GAN as the generative model. This synthetic experiment is illustrated in Figure 10. Specifically, we created a “real” distribution that is a Gaussian mixture distribution biased towards one mode. We derive a closed-form fair classifier based on this distribution, detailed in the purple block in Figure 10. When a GAN is trained to approximate this distribution, its outputs are biased towards the over-represented mode, as shown by Figure 10(b). Figure 10(c) illustrates the controllability experiment: using the fair classifier as control, PromptGen learns a distribution concentrated in a specific region of the output space. Figure 10(d) illustrates the de-biasing experiment: using the moment constraint (Section 3.1 and Section 4.2), PromptGen upweights the under-represented regions of the output space.

B Method Details and Derivations

B.1 Controllability as EBM

In Section 3, we defined a control \mathcal{C} as M independent properties $\{\mathbf{y}_1, \dots, \mathbf{y}_M\}$. For example, \mathbf{y}_1 can be a text description, and \mathbf{y}_2 can be an attribute. In this part, we first elaborate on why controllability can be formed as EBMs in Eq. (1). We then provide concrete examples of the distributions derived from different energy functions defined in Section 3.1.

We denote the image prior as $p_{\mathbf{x}}(\mathbf{x})$, the only distribution that can be estimated from data when labels for the control are not provided during generative model pre-training. Given the control $\mathcal{C} = \{\mathbf{y}_1, \dots, \mathbf{y}_M\}$, we resort to Bayes’ theorem to rewrite the conditional distribution $p(\mathbf{x}|\mathcal{C})$ as

$$p(\mathbf{x}|\mathcal{C}) \propto p_{\mathbf{x}}(\mathbf{x})p(\mathcal{C}|\mathbf{x}) \quad (13)$$

$$= p_{\mathbf{x}}(\mathbf{x})p(\mathbf{y}_1|\mathbf{x}) \prod_{i=2}^M p(\mathbf{y}_i|\mathbf{x}, \mathbf{y}_{<i}) \quad (14)$$

$$= p_{\mathbf{x}}(\mathbf{x}) \prod_{i=1}^M p(\mathbf{y}_i|\mathbf{x}) \quad (\text{independence assumption}). \quad (15)$$

For Eq. (15) to be well-defined, we need to define $p(\mathbf{y}_i|\mathbf{x})$ for each \mathbf{y}_i . In order to incorporate the knowledge of arbitrary off-the-self models (besides image classifiers), we define each $p(\mathbf{y}_i|\mathbf{x})$ as

$$p(\mathbf{y}_i|\mathbf{x}) = \frac{\exp(-\lambda_i E_i(\mathbf{x}, \mathbf{y}_i))}{Z_i}, \quad Z_i = \int_{\mathbf{y}'_i} \exp(-\lambda_i E_i(\mathbf{x}, \mathbf{y}'_i)) d\mathbf{y}'_i. \quad (16)$$

Eq. (16) says that $p(\mathbf{y}_i|\mathbf{x})$ is proportional to the exponential of an energy function $E_i(\mathbf{x}, \mathbf{y}_i)$, where \mathbf{y}_i with lower energy has higher density or mass. Note that Eq. (16) defines a distribution over all possible values of \mathbf{y}_i instead of all possible values of \mathbf{x} . Combining Eq. (15) and Eq. (16), we have

$$p(\mathbf{x}|\mathcal{C}) = \frac{p_{\mathbf{x}}(\mathbf{x})e^{-E_{\mathcal{C}}(\mathbf{x})}}{Z_X}, \quad E_{\mathcal{C}}(\mathbf{x}) = \sum_{i=1}^M \lambda_i E_i(\mathbf{x}, \mathbf{y}_i), \quad Z_X = \int_{\mathbf{x}'} p_{\mathbf{x}}(\mathbf{x}')e^{-E_{\mathcal{C}}(\mathbf{x}')} d\mathbf{x}', \quad (17)$$

which is the same equation as Eq. (1).

In the following, we use Eq. (16) to derive the distributions from the classifier energy, CLIP energy, and inverse graphics energy defined in Section 3.1.

Classifier energy: Given a classifier $P(\cdot|\mathbf{x})$ and the target class a , we define the classifier energy as $E_{\text{classifier}}(\mathbf{x}, a) = -\log P(a|\mathbf{x})$. Using Eq. (16), we arrive at:

$$p_{\text{classifier}}(a|\mathbf{x}) = \frac{\exp(\lambda_{\text{classifier}} \log P(a|\mathbf{x}))}{\sum_{a'} \exp(\lambda_{\text{classifier}} \log P(a'|\mathbf{x}))} = \frac{P(a|\mathbf{x})^{\lambda_{\text{classifier}}}}{\sum_{a'} P(a'|\mathbf{x})^{\lambda_{\text{classifier}}}}, \quad (18)$$

which is equivalent to a temperature-adjusted distribution of the original classifier.

CLIP energy: Using Eq. (16), the CLIP energy in Eq. (3) is equivalent to

$$p_{\text{CLIP}}(\mathbf{t}|\mathbf{x}) \propto \exp\left(-\frac{\lambda_{\text{CLIP}}}{L} \sum_{l=1}^L \left(1 - \cos \langle \text{CLIP}_{\text{img}}(\text{DiffAug}_l(\mathbf{x})), \text{CLIP}_{\text{text}}(\mathbf{t}) \rangle\right)\right). \quad (19)$$

Inverse graphics energy: Using Eq. (16), the inverse graphics energy in Eq. (4) is equivalent to

$$p_{\text{inv-graphics}}(\boldsymbol{\rho}|\mathbf{x}) \propto \exp(-\lambda_{\text{inv-graphics}} d(f_{\mathcal{X} \rightarrow \mathcal{P}}(\mathbf{x}), \boldsymbol{\rho})^2). \quad (20)$$

If the geodesic distance $d\langle \cdot, \cdot \rangle$ is the Euclidean distance, then $p_{\text{inv-graphics}}(\boldsymbol{\rho}|\mathbf{x})$ is a Gaussian distribution whose mean is $\boldsymbol{\rho}$ and whose variance depends on the hyperparameter $\lambda_{\text{inv-graphics}}$; if the geodesic distance $d\langle \cdot, \cdot \rangle$ is the spherical distance (e.g., the distance between pose parameters defined on a unit sphere), then $p_{\text{inv-graphics}}(\boldsymbol{\rho}|\mathbf{x})$ is a vMF distribution.

B.2 Equivalence between Image-Space EBM and Latent-Space EBM

Proposition 1. *Define $\mathbf{x} \sim p_{\mathbf{x}}(\mathbf{x})$ as $\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$, $\mathbf{x} = G(\mathbf{z})$ and $p(\mathbf{x}|\mathcal{C})$ as $\mathbf{z} \sim p(\mathbf{z}|\mathcal{C})$, $\mathbf{x} = G(\mathbf{z})$, where $p(\mathbf{z}|\mathcal{C})$ is defined as the following EBM:*

$$p(\mathbf{z}|\mathcal{C}) = \frac{p_{\mathbf{z}}(\mathbf{z})e^{-E_{\mathcal{C}}(G(\mathbf{z}))}}{Z}, \quad E_{\mathcal{C}}(\mathbf{x}) = \sum_{i=1}^M \lambda_i E_i(\mathbf{x}, \mathbf{y}_i), \quad Z = \int_{\mathbf{z}'} p_{\mathbf{z}}(\mathbf{z}')e^{-E_{\mathcal{C}}(G(\mathbf{z}'))} d\mathbf{z}'. \quad (21)$$

We have

$$p(\mathbf{x}|\mathcal{C}) = \frac{p_{\mathbf{x}}(\mathbf{x})e^{-E_{\mathcal{C}}(\mathbf{x})}}{Z_X}, \quad Z_X = \int_{\mathbf{x}'} p_{\mathbf{x}}(\mathbf{x}')e^{-E_{\mathcal{C}}(\mathbf{x}')} d\mathbf{x}'. \quad (22)$$

In spirit, our proof follows the proof in [50], which follows [19]. The difference between our proof and that in [50] is that we derive $p(\mathbf{x}|\mathcal{C})$ from $p(\mathbf{z}|\mathcal{C})$ while they derived $p(\mathbf{z}|\mathcal{C})$ from $p(\mathbf{x}|\mathcal{C})$.

Proof. Based on Lemma 1 in [19] and Lemma 1 in [50], $p(\mathbf{z}|\mathcal{C})$ is equivalent to rejection sampling with proposal distribution $p_{\mathbf{z}}(\mathbf{z})$ and acceptance probability

$$r(\mathbf{z}) = \frac{e^{-E_{\mathcal{C}}(G(\mathbf{z}))}}{M_{\mathcal{C}} \cdot Z}, \quad \text{where } \forall \mathbf{z}, M_{\mathcal{C}} > \frac{e^{-E_{\mathcal{C}}(G(\mathbf{z}))}}{Z}. \quad (23)$$

Since $p(\mathbf{x}|\mathcal{C})$ is defined as $\mathbf{z} \sim p(\mathbf{z}|\mathcal{C})$, $\mathbf{x} = G(\mathbf{z})$, $p(\mathbf{x}|\mathcal{C})$ is equivalent to rejection sampling with proposal distribution $p_{\mathbf{z}}(\mathbf{z})$, $\mathbf{x} = G(\mathbf{z})$ and acceptance probability

$$r(\mathbf{x}) = \frac{e^{-E_{\mathcal{C}}(\mathbf{x})}}{M_{\mathcal{C}} \cdot Z}. \quad (24)$$

Note that the above proposal distribution $\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$, $\mathbf{x} = G(\mathbf{z})$ is the same as $p_{\mathbf{x}}(\mathbf{x})$ by definition. Based on Lemma 1 in [19] and Lemma 1 in [50], we arrive at:

$$p(\mathbf{x}|\mathcal{C}) = \frac{p_{\mathbf{x}}(\mathbf{x})r(\mathbf{x})}{\mathbb{E}_{\mathbf{x}' \sim p_{\mathbf{x}}(\mathbf{x}')} [r(\mathbf{x}')] } \quad (25)$$

$$= \frac{p_{\mathbf{x}}(\mathbf{x})e^{-E_{\mathcal{C}}(\mathbf{x})}/(M_{\mathcal{C}} \cdot Z)}{\mathbb{E}_{\mathbf{x}' \sim p_{\mathbf{x}}(\mathbf{x}')} [e^{-E_{\mathcal{C}}(\mathbf{x}')}/(M_{\mathcal{C}} \cdot Z)]} \quad (26)$$

$$= \frac{p_{\mathbf{x}}(\mathbf{x})e^{-E_{\mathcal{C}}(\mathbf{x})}}{\mathbb{E}_{\mathbf{x}' \sim p_{\mathbf{x}}(\mathbf{x}')} [e^{-E_{\mathcal{C}}(\mathbf{x}')}] } \quad (27)$$

$$= \frac{p_{\mathbf{x}}(\mathbf{x})e^{-E_{\mathcal{C}}(\mathbf{x})}}{Z_X}, \quad Z_X = \int_{\mathbf{x}'} p_{\mathbf{x}}(\mathbf{x}')e^{-E_{\mathcal{C}}(\mathbf{x}')} d\mathbf{x}'. \quad (28)$$

□

B.3 Derivation of Eq. (9): Approximating EBM with INN

The full derivation of Eq. (9) is given by:

$$\begin{aligned} & \mathbb{D}_{\text{KL}}(p_{\theta}(\mathbf{z}) \| p(\mathbf{z}|\mathcal{C})) \\ &= \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{z})}{p(\mathbf{z}|\mathcal{C})} \right] \\ &= \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\mathbf{z}), \mathbf{x}=G(\mathbf{z})} \left[\log \frac{p_{\theta}(\mathbf{z})}{p_{\mathbf{z}}(\mathbf{z})e^{-E_{\mathcal{C}}(\mathbf{x})}/Z} \right] \\ &= \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{z}=f_{\theta}(\boldsymbol{\epsilon}), \mathbf{x}=G(\mathbf{z})} \left[\log \mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \mathbf{I}) - \log \left| \det \left(\frac{\partial f_{\theta}}{\partial \boldsymbol{\epsilon}} \right) \right| \right. \\ & \quad \left. - \log p_{\mathbf{z}}(\mathbf{z}) + E_{\mathcal{C}}(\mathbf{x}) + \log Z \right] \quad (29) \\ &= \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{z}=f_{\theta}(\boldsymbol{\epsilon}), \mathbf{x}=G(\mathbf{z})} \left[- \log \left| \det \left(\frac{\partial f_{\theta}}{\partial \boldsymbol{\epsilon}} \right) \right| - \log p_{\mathbf{z}}(\mathbf{z}) \right. \\ & \quad \left. + E_{\mathcal{C}}(\mathbf{x}) \right] - \mathbb{H}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} + \log Z. \end{aligned}$$

B.4 Derivation of Eq. (11): Extension to Generative Models with a Class-Embedding Space

The full derivation of Eq. (11) is given by:

$$\begin{aligned} & \mathbb{D}_{\text{KL}}(p_{\theta}(\mathbf{z}, \mathbf{y}) \| p(\mathbf{z}, \mathbf{y}|\mathcal{C})) \\ &= \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim p_{\theta}(\mathbf{z}, \mathbf{y})} \left[\log \frac{p_{\theta}(\mathbf{z}, \mathbf{y})}{p(\mathbf{z}, \mathbf{y}|\mathcal{C})} \right] \\ &= \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim p_{\theta}(\mathbf{z}, \mathbf{y}), \mathbf{x}=G(\mathbf{z}, \mathbf{y})} \left[\log \frac{p_{\theta}(\mathbf{z}, \mathbf{y})}{p_{\mathbf{z}, \mathbf{y}}(\mathbf{z}, \mathbf{y})e^{-E_{\mathcal{C}}(\mathbf{x})}/Z} \right] \\ &= \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim p_{\theta}(\mathbf{z}, \mathbf{y}), \mathbf{x}=G(\mathbf{z}, \mathbf{y})} \left[\log \frac{p_{\theta}(\mathbf{z}, \mathbf{y})}{p_{\mathbf{z}}(\mathbf{z})p_{\mathbf{y}}(\mathbf{y})e^{-E_{\mathcal{C}}(\mathbf{x})}/Z} \right] \quad (\mathbf{z} \text{ and } \mathbf{y} \text{ are independent}) \\ &= \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I}), \mathbf{z}=f_{\theta}(\boldsymbol{\epsilon}), \mathbf{y}=h_{\theta}(\boldsymbol{\xi}), \mathbf{x}=G(\mathbf{z}, \mathbf{y})} \left[\log \mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \mathbf{I}) - \log \left| \det \left(\frac{\partial f_{\theta}}{\partial \boldsymbol{\epsilon}} \right) \right| \right. \\ & \quad \left. + \log \mathcal{N}(\boldsymbol{\xi}|\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I}) - \log \left| \det \left(\frac{\partial h_{\theta}}{\partial \boldsymbol{\xi}} \right) \right| - \log p_{\mathbf{z}}(\mathbf{z}) - \log p_{\mathbf{y}}(\mathbf{y}) + E_{\mathcal{C}}(\mathbf{x}) + \log Z \right] \\ &= \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \boldsymbol{\xi} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I}), \mathbf{z}=f_{\theta}(\boldsymbol{\epsilon}), \mathbf{y}=h_{\theta}(\boldsymbol{\xi}), \mathbf{x}=G(\mathbf{z}, \mathbf{y})} \left[- \log \left| \det \left(\frac{\partial f_{\theta}}{\partial \boldsymbol{\epsilon}} \right) \right| - \log \left| \det \left(\frac{\partial h_{\theta}}{\partial \boldsymbol{\xi}} \right) \right| \right. \\ & \quad \left. - \log p_{\mathbf{z}}(\mathbf{z}) - \log p_{\mathbf{y}}(\mathbf{y}) + E_{\mathcal{C}}(\mathbf{x}) \right] - \mathbb{H}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} - \mathbb{H}_{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})} + \log Z. \end{aligned} \quad (30)$$

Algorithm 3: Extension of Algorithm 2 to Conditional INN

while *not converged* **do**

1. Sample $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\rho \sim p_\rho(\rho)$
 2. Map ϵ to latent code $\mathbf{z} = f_\theta(\epsilon, \rho)$
 3. Map \mathbf{z} to an image $\mathbf{x} = G(\mathbf{z})$
 4. Optimize θ with gradient $\nabla_\theta \left(-\log \left| \det \left(\frac{\partial f_\theta}{\partial \epsilon} \right) \right| - \log p_z(\mathbf{z}) + E_{\mathcal{C}_\rho}(\mathbf{x}) \right)$
-

B.5 Extension to Conditional INN

Algorithm 3 illustrates how we extend the training algorithm (Algorithm 2) to conditional INN, which helps generalize to controls specified by continuous values ρ . The training objective becomes

$$\arg \min_{\theta} \mathbb{E}_{\rho \sim p_\rho(\rho)} [\mathbb{D}_{\text{KL}}(p_\theta(\mathbf{z}|\rho) \| p(\mathbf{z}|\mathcal{C}_\rho))], \quad (31)$$

where \mathcal{C}_ρ means that the control is specified by value ρ . The derivation of Algorithm 3 is given by:

$$\begin{aligned} & \mathbb{E}_{\rho \sim p_\rho(\rho)} [\mathbb{D}_{\text{KL}}(p_\theta(\mathbf{z}|\rho) \| p(\mathbf{z}|\mathcal{C}_\rho))] \\ &= \mathbb{E}_{\rho \sim p_\rho(\rho), \mathbf{z} \sim p_\theta(\mathbf{z}|\rho)} \left[\log \frac{p_\theta(\mathbf{z}|\rho)}{p(\mathbf{z}|\mathcal{C}_\rho)} \right] \\ &= \mathbb{E}_{\rho \sim p_\rho(\rho), \mathbf{z} \sim p_\theta(\mathbf{z}|\rho), \mathbf{x} = G(\mathbf{z})} \left[\log \frac{p_\theta(\mathbf{z}|\rho)}{p_z(\mathbf{z}) e^{-E_{\mathcal{C}_\rho}(\mathbf{x})} / Z} \right] \\ &= \mathbb{E}_{\rho \sim p_\rho(\rho), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{z} = f_\theta(\epsilon, \rho), \mathbf{x} = G(\mathbf{z})} \left[\log \mathcal{N}(\epsilon | \mathbf{0}, \mathbf{I}) - \log \left| \det \left(\frac{\partial f_\theta}{\partial \epsilon} \right) \right| \right. \\ & \quad \left. - \log p_z(\mathbf{z}) + E_{\mathcal{C}_\rho}(\mathbf{x}) + \log Z \right] \\ &= \mathbb{E}_{\rho \sim p_\rho(\rho), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{z} = f_\theta(\epsilon, \rho), \mathbf{x} = G(\mathbf{z})} \left[-\log \left| \det \left(\frac{\partial f_\theta}{\partial \epsilon} \right) \right| - \log p_z(\mathbf{z}) \right. \\ & \quad \left. + E_{\mathcal{C}_\rho}(\mathbf{x}) \right] - \mathbb{H}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} + \log Z. \end{aligned} \quad (32)$$

B.6 Derivation for Latent-Space Moment Constraint

Given a mapping $\gamma : \mathcal{X} \rightarrow \mathbb{R}^K$, moment constraint [7, 36] defines the target distribution $p(\mathbf{x}|\mathcal{C})$ as

$$p(\mathbf{x}|\mathcal{C}) = \arg \min_{p(\mathbf{x}|\mathcal{C})} \mathbb{D}_{\text{KL}}(p(\mathbf{x}|\mathcal{C}) \| p_x(\mathbf{x})), \quad \text{s.t. } \underbrace{\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathcal{C})} [\gamma(\mathbf{x})]}_{\text{Moment constraint}} = \boldsymbol{\mu}, \quad (33)$$

Deviation from the pre-trained distribution

where $\boldsymbol{\mu}$ is the user-specified constraint. Examples are provided in the main text, and we omit them here for brevity. [7] showed that Eq. (33) can be approximated to an arbitrary precision by

$$p(\mathbf{x}|\mathcal{C}) \propto p_x(\mathbf{x}) \exp(\hat{\boldsymbol{\beta}}^\top \gamma(\mathbf{x})), \quad (34)$$

where $\hat{\boldsymbol{\beta}}$ needs to be computed. [36] estimates $\hat{\boldsymbol{\beta}}$ by solving the following regression problem:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \mathbb{E}_{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)} \stackrel{\text{i.i.d.}}{\sim} p_x(\mathbf{x})} \left\| \frac{\sum_{j=1}^N \exp(\boldsymbol{\beta}^\top \gamma(\mathbf{x}^{(j)})) \gamma(\mathbf{x}^{(j)})}{\sum_{j'=1}^N \exp(\boldsymbol{\beta}^\top \gamma(\mathbf{x}^{(j')}))} - \boldsymbol{\mu} \right\|_2^2. \quad (35)$$

In [36], sampling $\mathbf{x} \sim p_x(\mathbf{x})$ is straightforward since they focus on autoregressive language models. In the context of latent-variable vision generative models, where $\mathbf{x} \sim p_x(\mathbf{x})$ is implicitly defined as $\mathbf{z} \sim p_z(\mathbf{z})$, $\mathbf{x} = G(\mathbf{z})$, Eq. (35) is equivalent to

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \mathbb{E}_{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(N)} \stackrel{\text{i.i.d.}}{\sim} p_z(\mathbf{z}), \mathbf{x}^{(j)} = G(\mathbf{z}^{(j)})} \left\| \frac{\sum_{j=1}^N \exp(\boldsymbol{\beta}^\top \gamma(\mathbf{x}^{(j)})) \gamma(\mathbf{x}^{(j)})}{\sum_{j'=1}^N \exp(\boldsymbol{\beta}^\top \gamma(\mathbf{x}^{(j')}))} - \boldsymbol{\mu} \right\|_2^2. \quad (36)$$

Finally, based on Proposition 1, we can generalize the moment constraint to the latent space as

$$p(\mathbf{z}|\mathcal{C}) = \frac{p_{\mathbf{z}}(\mathbf{z}) \exp\left(\hat{\boldsymbol{\beta}}^\top \boldsymbol{\gamma}(G(\mathbf{z}))\right)}{Z}, \quad Z = \int_{\mathbf{z}'} p_{\mathbf{z}}(\mathbf{z}') \exp\left(\hat{\boldsymbol{\beta}}^\top \boldsymbol{\gamma}(G(\mathbf{z}'))\right) d\mathbf{z}'. \quad (37)$$

Note that Eq. (36) and Eq. (37) are verbatim copies of Eq. (7) and Eq. (6), respectively.

Algorithm 4: Extension of Algorithm 3 with ID Energy (Section 4.3)

while not converged do

1. Sample $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\boldsymbol{\rho} \sim p_{\boldsymbol{\rho}}(\boldsymbol{\rho})$
2. Map $\boldsymbol{\epsilon}$ to latent codes $\mathbf{z} = f_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}, \boldsymbol{\rho})$ and $\mathbf{z}_0 = f_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}, \boldsymbol{\rho}_0)$
3. Map \mathbf{z} to an image $\mathbf{x} = G(\mathbf{z})$ and \mathbf{z}_0 to an image $\mathbf{x}_0 = G(\mathbf{z}_0)$
4. Optimize $\boldsymbol{\theta}$ with gradient

$$\nabla_{\boldsymbol{\theta}} \left(-\log \left| \det \left(\frac{\partial f_{\boldsymbol{\theta}}}{\partial \boldsymbol{\epsilon}} \right) \right| - \log p_{\mathbf{z}}(\mathbf{z}) + E_{C_{\rho}}(\mathbf{x}) + \lambda_{\text{ID}} E_{\text{ID}}(\mathbf{x}_0, \mathbf{x}) \right)$$

B.7 Inverse Graphics Control with Identity Energy

This subsection provides more detail on how we use an inverse graphics model, DECA [16], to control the pose of faces generated by StyleGAN2 trained on FFHQ 1024², as introduced in Section 4.3. Recall that we use the conditional INN extension, detailed in Algorithm 3, for this experiment. To enable generating different poses of the same identity, we add additional identity energy using the IR-SE50 model [8]. The training algorithm is detailed in Algorithm 4. Specifically, we generate a canonical latent code $\mathbf{z}_0 = f_{\boldsymbol{\theta}}(\boldsymbol{\epsilon}, \boldsymbol{\rho}_0)$ for each $\boldsymbol{\epsilon}$, where $\boldsymbol{\rho}_0$ is the canonical pose. The latent code \mathbf{z} and the canonical latent code \mathbf{z}_0 are mapped to the image \mathbf{x} and the canonical image \mathbf{x}_0 . The identity energy in Eq. (12), copied below, encourages the embeddings of the two images to be similar:

$$E_{\text{ID}}(\mathbf{x}_0, \mathbf{x}) = 1 - \cos \langle R(\mathbf{x}_0), R(\mathbf{x}) \rangle, \quad \mathbf{x}_0 = G(\mathbf{z}_0), \mathbf{x} = G(\mathbf{z}). \quad (38)$$

B.8 Experimental Details

Our INN architecture contains 8 blocks. Each block consists of a soft permutation of channels [2], an affine coupling layer [12], and an ActNorm layer [37]. In the affine coupling layer, we model the sub-network as an MLP with one hidden layer, where the hidden dimension is 256 and the non-linear activation is LeakyReLU_{0.1}.

For the prior distribution $p_{\boldsymbol{\rho}}(\boldsymbol{\rho})$ of the pose parameter $\boldsymbol{\rho}$ in the inverse graphics experiments (Section 4.3), we sample the x -axis and y -axis rotations relative to a canonical pose $\boldsymbol{\rho}_0$, whose x -axis and y -axis rotations are 0.3 and 0, respectively. The relative pose’s x -axis and y -axis rotations are uniformly sampled from $[-(1/9)\pi, (1/9)\pi]$.

Each experiment was run on an NVIDIA RTX A4000 GPU (with 16G memory). Our code implementation is based on the PyTorch framework. Our code can be found at <https://github.com/ChenWu98/Generative-Visual-Prompt>.

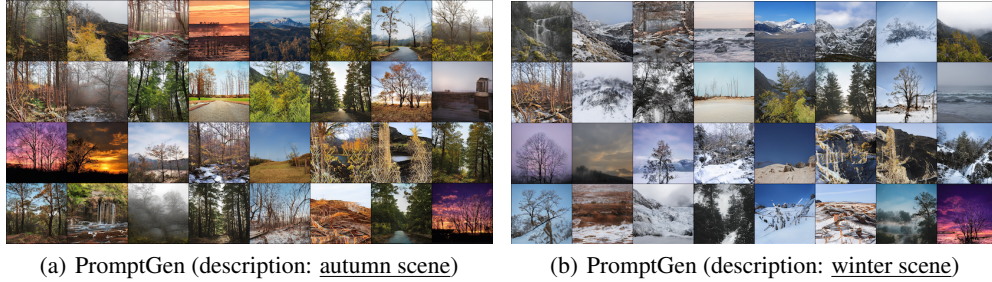


Figure 11: Additional results for image synthesis based on text description, guided by the CLIP model. As the backbone, we use StyleGAN2 trained on Landscape-HQ with truncation $\psi = 0.7$. All images are 256^2 and resized for visualization.

Table 5: Quantitative results for de-biasing categorical attributes (Figure 6). See details in Appendix D. PromptGen de-biases StyleGAN2 in terms of race, age, and gender.

	FFHQ		MetFaces		
	$\mathbb{D}_{KL}^{\text{face}} \downarrow$	$\mathbb{D}_{KL}^{\text{age}} \downarrow$	$\mathbb{D}_{KL}^{\text{race}} \downarrow$	$\mathbb{D}_{KL}^{\text{age}} \downarrow$	$\mathbb{D}_{KL}^{\text{gender}} \downarrow$
StyleGAN2	0.860	0.597	1.624	0.546	0.019
PromptGen (ours; $\lambda = 1$)	0.286	0.357	0.687	0.397	0.000
PromptGen (ours; $\lambda = 2$)	0.099	0.172	0.189	0.247	0.005

C Additional Results for Image Synthesis based on Text Description

This section presents additional results of PromptGen for image synthesis from text, using the CLIP model as control. As the pre-trained generative model, we explore StyleGAN2 trained on AFHQ-Cats [6], LSUN-Churches [81], FFHQ [34], and Landscape-HQ [68]. Figure 11 and Figure 12 present more results of PromptGen with diverse text descriptions.

D Additional Results for De-Biasing Generative Models

Table 5 provides quantitative results of de-biasing StyleGAN2 across categorical attributes (Figure 6). Specifically, we used the pre-trained classifier provided by FairFace [31] to classify the attributes of the generated images. We then report the KL divergence between the attribute distribution of the generated images and the uniform distribution. Figure 13 and Figure 14 visualize the de-biasing results. We report the KL divergence between the generated distribution and the uniform distribution. Interestingly, on MetFaces, de-biasing the race results in more sculptures, and we postulate that the reason is that almost all paintings and sketches in MetFaces are for white individuals.



Figure 12: Additional results for image synthesis based on text description, guided by the CLIP model. As the pre-trained generative model, we use StyleGAN2 trained on FFHQ, AFHQ-Cats and LSUN-Churches datasets. The captions are the text descriptions given to the CLIP model.



(a) Pre-trained StyleGAN2



(b) PromptGen (race de-biasing; $\lambda = 1$)



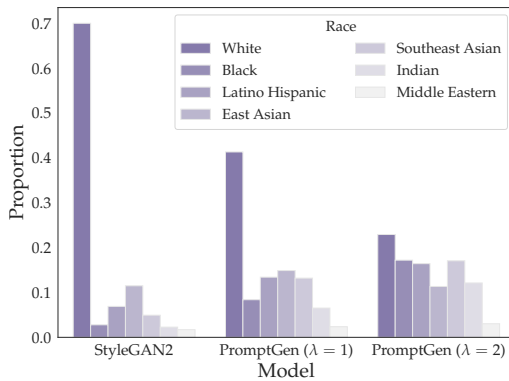
(c) PromptGen (race de-biasing; $\lambda = 2$)



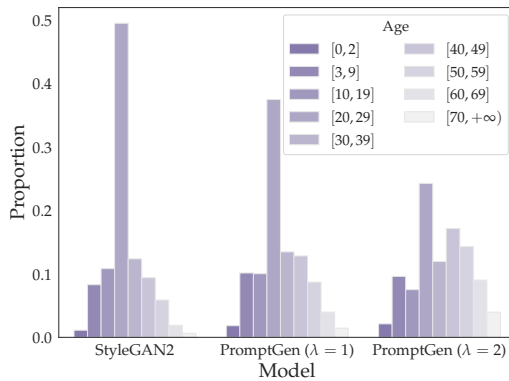
(d) PromptGen (age de-biasing; $\lambda = 1$)



(e) PromptGen (age de-biasing; $\lambda = 2$)



(f) Race distribution



(g) Age distribution

Figure 13: Using our moment constraint, PromptGen de-biases the racial and age distributions of StyleGAN2 trained on FFHQ with truncation $\psi = 0.7$. All synthesized images are 1024^2 in resolution and resized for visualization. We fixed the random seed for PromptGen, so please zoom in to see detailed differences between images.



(a) Pre-trained StyleGAN2



(b) PromptGen (race de-biasing; $\lambda = 1$)



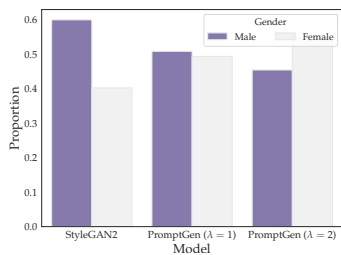
(c) PromptGen (race de-biasing; $\lambda = 2$)



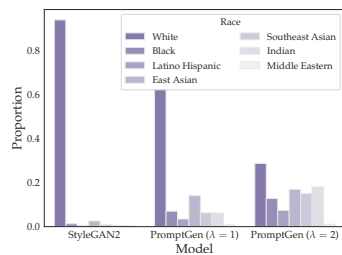
(d) PromptGen (age de-biasing; $\lambda = 2$)



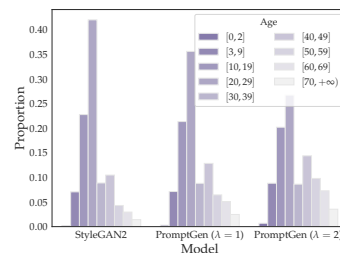
(e) PromptGen (gender de-biasing; $\lambda = 2$)



(f) Gender distribution



(g) Race distribution



(h) Age distribution

Figure 14: Using the moment constraint, PromptGen de-biases the racial, age, and gender distributions of StyleGAN2 trained on MetFaces with truncation $\psi = 0.7$. All images are 1024^2 and resized for visualization. Interestingly, de-biasing the race results in more sculptures, and we postulate that the reason is that almost all paintings and sketches in MetFaces are for white individuals. We fixed the random seed for PromptGen, so we recommend zooming in to see differences between images.

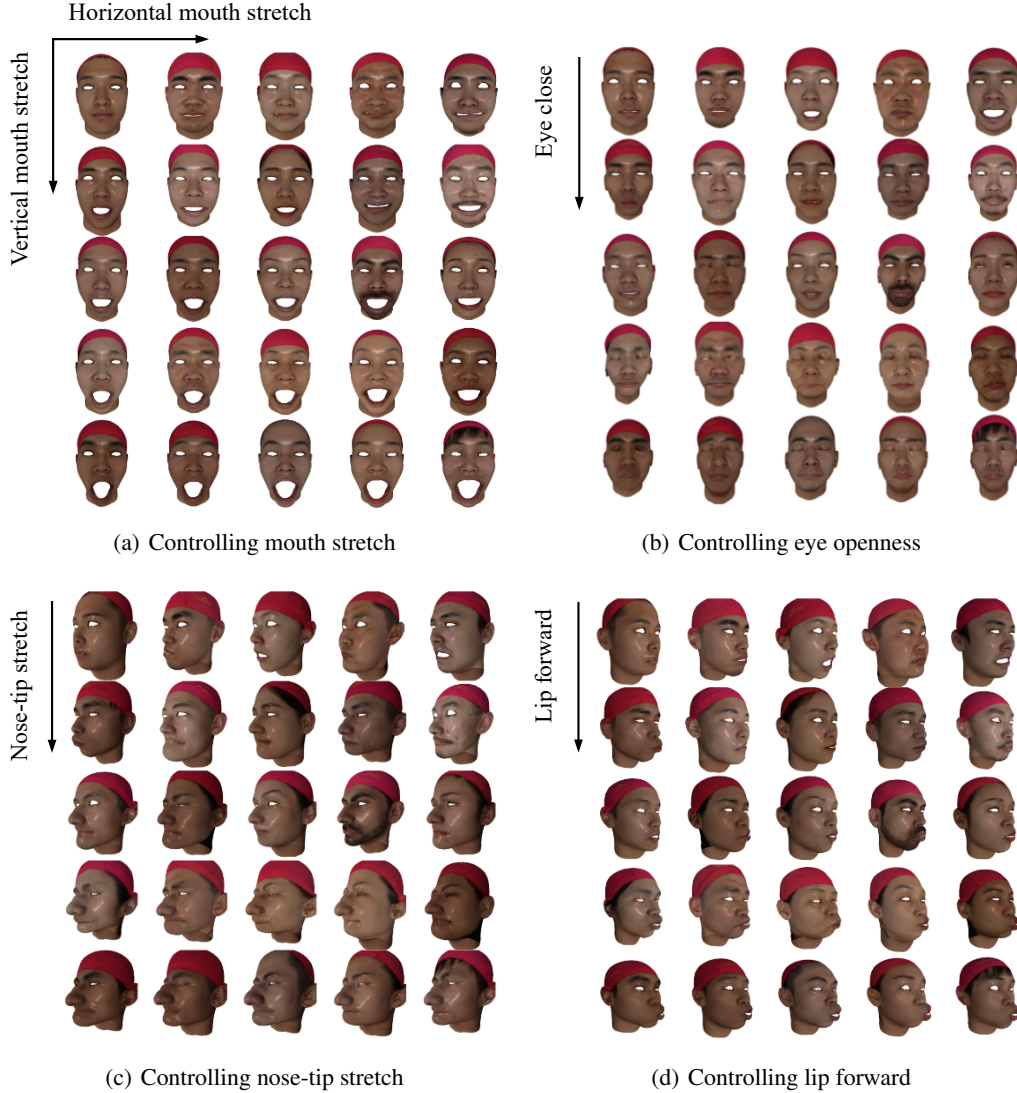


Figure 15: PromptGen with intra-sample signed-distance constraints that generate 3D face models with varying distances in the mouth, eye openness, nose-tip stretch, and lip forward.

E Experiments on Generative 3D Face Models

3D face modeling has been an active area of research that has recently gained major interest due to applications in virtual humans, deep faces, and digital actors. Existing 3D deep learning generative models build 3D compact representations of shape and appearance capable of modeling non-linearity (e.g., scatter effects, specularities) that is necessary for generating photo-realistic faces. Providing a generative model with the capability of generating 3D faces with a particular geometry (e.g., a specific distance between eyes or nose length) is particularly useful for technical artists when creating new characters. However, this problem remains unaddressed, partially due to the lack of publicly available 3D databases labeled with such constraints. We refer to this capability as the intra-sample constraint in 3D generative models, and this section shows how our PromptGen framework is able to achieve this by defining a new energy function based on signed distances.

We conducted the experiments on the FaceScape dataset [80], a large-scale 3D human face dataset consisting of 16,940 topologically uniformly registered 3D face meshes along with high-quality textures. The dataset contains 847 identities with 20 expressions per identity, and each mesh consists of 26,317 vertices of which 68 are 3D landmark vertices covering eyes, nose, mouth, jaw, and

eyebrows. With some abuse of notation, we denote a 3D face mesh as $\mathbf{x} \in \mathbb{R}^{3V}$ (recall that \mathbf{x} is used to represent an image in other parts of this paper), where $V = 26,317$ is the number of vertices. Given an index l of a vertex (e.g., one of the 68 landmark vertices), the 3D coordinate of this vertex is defined as $\mathbf{x}[l] \in \mathbb{R}^3$. We pre-train a generative mesh model G that maps each latent code $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to a 3D mesh \mathbf{x} . PromptGen allows us to obtain geometric control on the 3D meshes. Specifically, given two landmark vertex indices l_1 and l_2 , we define an intra-sample constraint that enforces the signed distance (along a given direction) between the two vertices to be s . To this effect, we define the signed-distance energy as

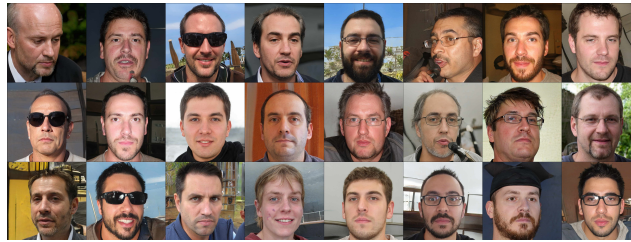
$$E_{\text{signed-distance}}(\mathbf{x}) = |d(\mathbf{v}_1, \mathbf{v}_2) - s|, \quad \mathbf{v}_1 = \mathbf{x}[l_1], \mathbf{v}_2 = \mathbf{x}[l_2] \quad (39)$$

where $d(\cdot, \cdot)$ is the signed-distance (along a given direction), and $|\cdot|$ is the absolute value.

Figure 15 presents several examples of 3D mesh control, in which random textures are applied to the generated meshes. It shows that PromptGen successfully controls the generative 3D mesh model in terms of the mouth stretch, eye openness, nose-tip stretch, and lip forward.

F Error Analysis

As we discussed in Section 5, the controlled distribution depends on (1) the pre-trained generative model’s coverage and (2) the off-the-shelf models used for the control. In this section, we provide some examples of PromptGen failing to generate satisfactory images using the CLIP energy. For the pre-trained generative model, we use StyleGAN2 trained on FFHQ, AFHQ-Wild, or LSUN-Churches.



(a) photo of a man without beard (FFHQ)



(b) photo of a person without makeup (FFHQ)

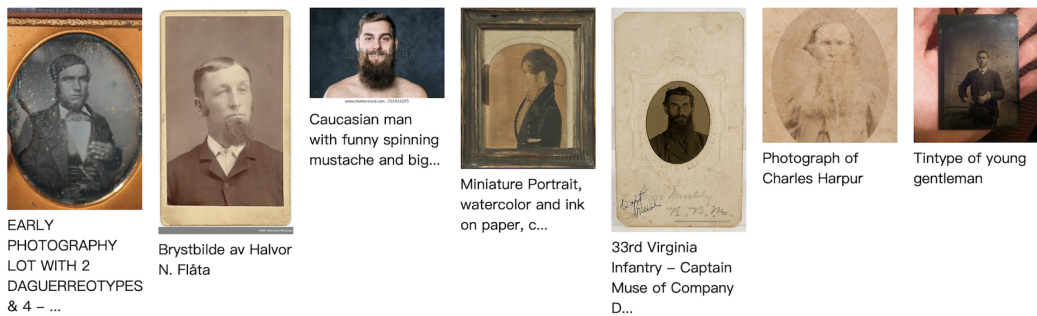


(c) photo of a church with three windows (LSUN-Churches)

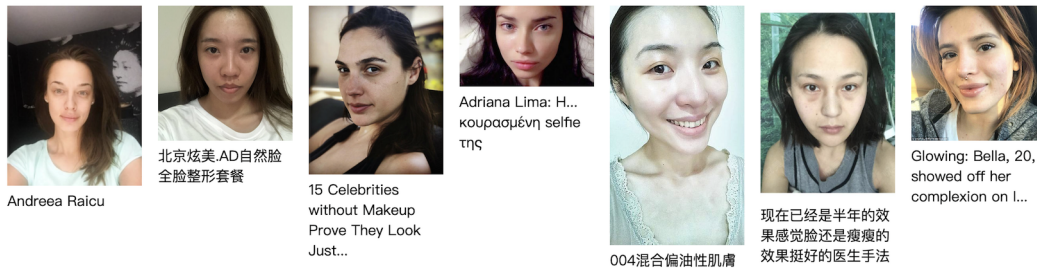
Figure 16: When the text description contains certain linguistic properties (e.g., negation, numerical reasoning), CLIP sometimes fails or shows the “reporting bias” that we discuss in Section 4.5. For a deeper understanding of these failures and biases, in Figure 17, we provide some image samples from the LAION-400M dataset [63] with CLIP retrieval, using the same text descriptions as this figure.

Our first observation is that the CLIP model sometimes fails in modeling linguistic negation. For instance, the text description photo of a man without beard results in a distribution of photos of men with beard (Figure 16(a)). Meanwhile, Figure 16(b) shows that CLIP is capable of modeling the negation of having makeup, but with the “reporting bias” discussed in Section 4.5. Moreover, CLIP seems to have difficulty in numerical reasoning, and gaining control over the count of specific objects in a scene tends to be unsuccessful. We showed this by specifying photo of a church with three windows, which did not result in the desired specification (Figure 16(c)).

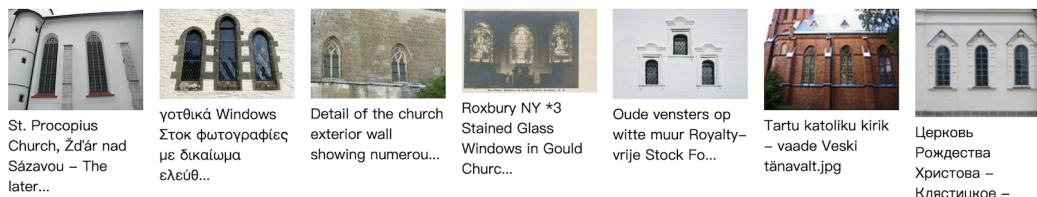
To gain a deeper understanding of the above failures and biases, in Figure 17, we provide some image samples from the LAION-400M dataset [63] with CLIP retrieval, using the same text descriptions as Figure 16. CLIP retrieval is based on the CLIP embedding similarity between the web images and text descriptions, while the original text below each individual image is not used. We observe an impressive consistency between CLIP retrieval and PromptGen: in both Figure 16(a) and Figure 17(a), most images have beard; in both Figure 16(b) and Figure 17(b), all images are female; in both Figure 16(c) and Figure 17(c), some images do not have exactly three windows. This consistency suggests that the failures and biases in Figure 16 should be mostly attributed to the CLIP model rather than to our PromptGen algorithm. We believe that our observation sheds light on the intricacy of contrastive multi-modal (vision-language) pre-training, which is worthy of being further investigated.



(a) photo of a man without beard (CLIP retrieval from LAION-400M)



(b) a photo of a person without makeup (CLIP retrieval from LAION-400M)



(c) photo of a church with three windows (CLIP retrieval from LAION-400M)

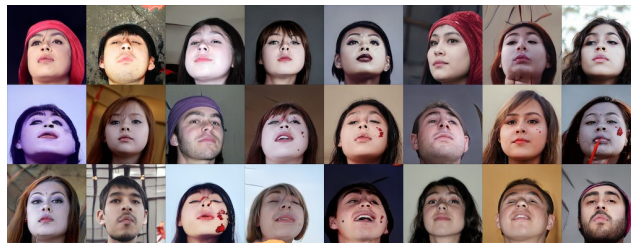
Figure 17: For a deeper understanding of the failures and biases illustrated in Figure 16, we provide some image samples from the LAION-400M dataset [63] with CLIP retrieval, using the same text descriptions as Figure 16. CLIP retrieval is based on the CLIP embedding similarity between the web images and text descriptions, while the original text below each individual image is not used. In Figure 17(a), most images have a beard; in Figure 17(b), all images are female; in Figure 17(c), some images do not have exactly three windows. These observations are consistent with those in Figure 16.

Another observation is that the control tends to fail when the text description requires sampling from low-density regions of the pre-trained generative model’s output space. In other words, the control usually fails if the pre-trained generative model does not cover the mode we are trying to gain control over. For example, images faithful to photo of a person yawning and photo of a baby with long hair are not commonly observed in the FFHQ dataset and, hence, these two text descriptions result in degeneration (Figure 18(a)) or weird images (Figure 18(b)). Another example is photo of an animal from the side, which is not commonly observed in the AFHQ-Wild dataset, and Figure 18(c) shows that the generated images fail to follow this description. Even when the control is successful (e.g., when a complex description is decomposed in Figure 8(b)), sampling from low-density regions results in limited diversity (e.g., the backgrounds in Figure 8(b) look similar to each other).

Finally, we also ran some failure controls using PPGM, showing that PromptGen and PPGM reveal similar failure cases of pre-trained generative models and CLIP. Results are shown in Figure 19.



(a) photo of a person yawning (FFHQ)



(b) photo of a baby with long hair (FFHQ)



(c) photo of an animal from the side (AFHQ-Wild)

Figure 18: When the pre-trained generative model fails in covering certain modes required by the text description, unsatisfactory outputs are produced. In this figure, we show several text descriptions that require sampling from low-density regions of the pre-trained generative model’s output space.



(a) photo of a man without beard (FFHQ; w/ PPGM) (b) photo of a person yawning (FFHQ; w/ PPGM)

Figure 19: Failure modes revealed by PromptGen (Figure 16 and Figure 18) also hold for PPGM.

G Societal Impact

With the improvements in generative models for images, DeepFake technology [75, 72] has become more accessible. Like any new technology, it is a double-edged sword, and it is crucial to research and comprehend the possible advantages and disadvantages of generative models for society.

On the positive side, we show that PromptGen can be used to de-bias pre-trained generative models and to reveal biases learned by text-image models (like CLIP), indicating that PromptGen might be a useful tool for fair AI if used appropriately. The efficient inference provided by PromptGen also helps reduce the computational expense, which has a positive impact on the environment. Better controllability, however, unavoidably makes it simpler to synthesize targeted pictures, which might have detrimental social effects in creating deceptive media (e.g., DeepFakes) or privacy leaks (e.g., identity-conditioned human face synthesis). To battle these cases, there are current technologies that can detect fake media effectively, and we expect companies and users to use these technologies to distinguish what is real from fake. We encourage practitioners to consider these risks when using PromptGen to develop systems.