

# RG-FLOW: A HIERARCHICAL AND EXPLAINABLE FLOW MODEL BASED ON RENORMALIZATION GROUP AND SPARSE PRIOR

**Hong-Ye Hu\*** **Dian Wu\*** **Yi-Zhuang You**

Department of Physics

UC San Diego

{hyhu, yzyou}@ucsd.edu

**Bruno Olshausen** **Yubei Chen**

Redwood Center, Berkeley AI Research

UC Berkeley

{baolshausen, yubeic}@berkeley.edu

## ABSTRACT

Flow-based generative models have become an important class of unsupervised learning approaches. In this work, we incorporate the key idea of renormalization group (RG) and sparse prior distribution to design a hierarchical flow-based generative model, called RG-Flow, which can separate information at different scales of images with disentangled representations at each scale. We demonstrate our method mainly on the CelebA dataset and show that the disentangled representations at different scales enable semantic manipulation and style mixing of the images. To visualize the latent representations, we introduce receptive fields for flow-based models and find that the receptive fields learned by RG-Flow are similar to those in convolutional neural networks. In addition, we replace the widely adopted Gaussian prior distribution by a sparse prior distribution to further enhance the disentanglement of representations. From a theoretical perspective, the proposed method has  $O(\log L)$  complexity for image inpainting compared to previous generative models with  $O(L^2)$  complexity.

## 1 INTRODUCTION

One of the most important unsupervised learning tasks is to learn the data distribution and build generative models. Over the past few years, various types of generative models have been proposed. Flow-based generative models are a particular family of generative models with tractable distributions (Dinh et al., 2017; Kingma & Dhariwal, 2018; Chen et al., 2018b; 2019; Behrmann et al., 2019; Hoogeboom et al., 2019; Brehmer & Cranmer, 2020; Rezende et al., 2020; Karami et al., 2019). Yet the latent variables are on equal footing and mixed globally. Here, we propose a new flow-based model, RG-Flow, which is inspired by the idea of *renormalization group* in statistical physics. RG-Flow imposes locality and hierarchical structure in bijective transformations. It allows us to access information at different scales in original images by latent variables at different locations, which offers better explainability. Combined with sparse priors (Olshausen & Field, 1996; 1997; Hyvärinen & Oja, 2000), we show that RG-Flow achieves hierarchical disentangled representations.

Renormalization group (RG) is a powerful tool to analyze statistical mechanics models and quantum field theories in physics (Kadanoff, 1966; Wilson, 1971). It progressively extracts more coarse-scale statistical features of the physical system and decimates irrelevant fine-grained statistics at each scale. Typically, the local transformations used in RG are designed by human physicists and they are not bijective. On the other hand, the flow-based models use cascaded invertible global transformations to progressively turn a complicated data distribution into Gaussian distribution. Here, we

\*H.-Y. Hu and D. Wu contribute equally to this work. Correspondence to: hyhu@ucsd.edu

would like to combine the key ideas from RG and flow-based models. The proposed RG-flow enables the machine to learn the optimal RG transformation from data, by constructing local invertible transformations and build a hierarchical generative model for the data distribution. Latent representations are introduced at different scales, which capture the statistical features at the corresponding scales. Together, the latent representations of all scales can be jointly inverted to generate the data. This method was recently proposed in the physics community as NeuralRG (Li & Wang, 2018; Hu et al., 2020).

Our main contributions are two-fold: First, RG-Flow can separate the signal statistics of different scales in the input distribution naturally, and represent information at each scale in its latent variables  $z$ . Those hierarchical latent variables live on a hyperbolic tree. Taking CelebA dataset (Liu et al., 2015) as an example, the network will not only find high-level representations, such as the gender factor and the emotion factor for human faces, but also mid-level and low-level representations. To visualize representations of different scales, we adopt the concept of *receptive field* from convolutional neural networks (CNN) (LeCun, 1988; LeCun et al., 1989) and visualize the hidden structures in RG-flow. In addition, since the statistics are separated into a hierarchical fashion, we show that the representations can be mixed at different scales. This achieves an effect similar to style mixing. Second, we introduce the *sparse prior distribution* for latent variables. We find the sparse prior distribution is helpful to further disentangle representations and make them more explainable. The widely adopted Gaussian prior is rotationally symmetric. As a result, each of the latent variables in a flow model usually does not have a clear semantic meaning. By using a sparse prior, we demonstrate the clear semantic meaning in the latent space.

## 2 RELATED WORK

Some flow-based generative models also possess multi-scale latent space (Dinh et al., 2017; Kingma & Dhariwal, 2018), and recently hierarchies of features have been utilized in Schirmer et al. (2020), where the top-level feature is shown to perform strongly in out-of-distribution (OOD) detection task. Yet, previous models do not impose hard locality constraint in the multi-scale structure. In Appendix C, the differences between globally connected multi-scale flows and RG-Flow are discussed, and we see that semantic, meaningful receptive fields do not show up in the globally connected cases. Recently, other more expressive bijective maps have been developed (Hoogeboom et al., 2019; Karami et al., 2019; Durkan et al., 2019), and those methods can be incorporated into the proposed structure to further improve the expressiveness of RG-Flow.

Some other classes of generative models rely on a separate inference model to obtain the latent representation. Examples include variational autoencoders (Kingma & Welling, 2014), adversarial autoencoders (Makhzani et al., 2015), InfoGAN (Chen et al., 2016), and BiGAN (Donahue et al., 2017; Dumoulin et al., 2017). Those techniques typically do not use hierarchical latent variables, and the inference of latent variables is approximate. Notably, recent advances suggest that having hierarchical latent variables may be beneficial (Vahdat & Kautz, 2020). In addition, the coarse-to-fine fashion of the generation process has also been discussed in other generative models, such as Laplacian pyramid of adversarial networks (Denton et al., 2015), and multi-scale autoregressive models (Reed et al., 2017).

*Disentangled representations* (Tenenbaum & Freeman, 2000; DiCarlo & Cox, 2007; Bengio et al., 2013) is another important aspect in understanding how a model generates images (Higgins et al., 2018). Especially, disentangled high-level representations have been discussed and improved from information theoretical principles (Cheung et al., 2015; Chen et al., 2016; 2018a; Higgins et al., 2017; Kipf et al., 2020; Kim & Mnih, 2018; Locatello et al., 2019; Ramesh et al., 2018). Apart from the high-level representations, the multi-scale structure also lies in the distribution of natural images. If a model can separate information of different scales, then its multi-scale representations can be used to perform other tasks, such as style transfer (Gatys et al., 2016; Zhu et al., 2017), face mixing (Karras et al., 2019; Gambardella et al., 2019; Karras et al., 2020), and texture synthesis (Bergmann et al., 2017; Jetchev et al., 2016; Gatys et al., 2015; Johnson et al., 2016; Ulyanov et al., 2016).

Typically, in flow-based generative models, Gaussian distribution is used as the prior for the latent space. Due to the rotational symmetry of Gaussian prior, an arbitrary rotation of the latent space would lead to the same likelihood. Sparse priors (Olshausen & Field, 1996; 1997; Hyvärinen & Oja, 2000) was proposed as an important tool for unsupervised learning and it leads to better explain-

ability in various domains (Ainsworth et al., 2018; Arora et al., 2018; Zhang et al., 2019). To break the symmetry of Gaussian prior and further improve the explainability, we introduce a sparse prior to flow-based models. Please refer to Figure 12 for a quick illustration on the difference between Gaussian prior and the sparse prior, where the sparse prior leads to better disentanglement.

*Renormalization group* (RG) has a broad impact ranging from particle physics to statistical physics. Apart from the analytical studies in field theories (Wilson, 1971; Fisher, 1998; Stanley, 1999), RG has also been useful in numerically simulating quantum states. The multi-scale entanglement renormalization ansatz (MERA) (Vidal, 2008; Evenbly & Vidal, 2014) implements the hierarchical structure of RG in tensor networks to represent quantum states. The exact holographic mapping (EHM) (Qi, 2013; Lee & Qi, 2016; You et al., 2016) further extends MERA to a bijective (unitary) flow between latent product states and visible entangled states. Recently, Li & Wang (2018); Hu et al. (2020) incorporates the MERA structure and deep neural networks to design a flow-based generative model that allows machine to learn the EHM from statistical physics and quantum field theory actions. In quantum machine learning, recent development of quantum convolutional neural networks also (Cong et al., 2019) utilize the MERA structure. The similarity between RG and deep learning has been discussed in several works (Bény, 2013; Mehta & Schwab, 2014; Bény & Osborne, 2015; Oprisa & Toth, 2017; Lin et al., 2017; Gan & Shu, 2017). The information theoretic objective that guides machine-learning RG transforms are proposed in recent works (Koch-Janusz & Ringel, 2018; Hu et al., 2020; Lenggenhager et al., 2020). The meaning of the emergent latent space has been related to quantum gravity (Swingle, 2012; Pastawski et al., 2015), which leads to the exciting development of machine learning holography (You et al., 2018; Hashimoto et al., 2018; Hashimoto, 2019; Akutagawa et al., 2020; Hashimoto et al., 2020).

### 3 METHODS

**Flow-based generative models.** Flow-based generative models are a family of generative models with tractable distributions, which allows efficient sampling and exact evaluation of the probability density (Dinh et al., 2015; 2017; Kingma & Dhariwal, 2018; Chen et al., 2019). The key idea is to build a bijective map  $G(\mathbf{z}) = \mathbf{x}$  between visible variables  $\mathbf{x}$  and latent variables  $\mathbf{z}$ . Visible variables  $\mathbf{x}$  are the data that we want to generate, which may follow a complicated probability distribution. And latent variables  $\mathbf{z}$  usually have simple distribution that can be easily sampled, for example the i.i.d. Gaussian distribution. In this way, the data can be efficiently generated by first sampling  $\mathbf{z}$  and mapping them to  $\mathbf{x}$  through  $\mathbf{x} = G(\mathbf{z})$ . In addition, we can get the probability associated with each data sample  $\mathbf{x}$ ,

$$\log p_X(\mathbf{x}) = \log p_Z(\mathbf{z}) - \log \left| \frac{\partial G(\mathbf{z})}{\partial \mathbf{z}} \right|. \quad (1)$$

The bijective map  $G(\mathbf{z}) = \mathbf{x}$  is usually composed as a series of bijectors,  $G(\mathbf{z}) = G_1 \circ G_2 \circ \dots \circ G_n(\mathbf{z})$ , such that each bijector layer  $G_i$  has a tractable Jacobian determinant and can be inverted efficiently. The two key ingredients in flow-based models are the design of the bijective map  $G$  and the choice of the prior distribution  $p_Z(\mathbf{z})$ .

**Structure of RG-Flow networks.** Much of the prior research has focused on designing more powerful bijective blocks for the generator  $G$  to improve its expressive power and to achieve better approximations of complicated probability distributions. Here, we focus on designing the architecture that arranges the bijective blocks in a hierarchical structure to separate features of different scales in the data and to disentangle latent representations.

Our design is motivated by the idea of RG in physics, which progressively separates the coarse-grained data statistics from fine-grained statistics by local transformations at different scales. Let  $\mathbf{x}$  be the visible variables, or the input image (level-0), denoted as  $\mathbf{x}^{(0)} \equiv \mathbf{x}$ . A step of the RG transformation extracts the coarse-grained information  $\mathbf{x}^{(1)}$  to send to the next layer (level-1), and splits out the rest of fine-grained information as auxiliary variables  $\mathbf{z}^{(0)}$ . The procedure can be described by the following recursive equation (at level- $h$  for example),

$$\mathbf{x}^{(h+1)}, \mathbf{z}^{(h)} = R_h(\mathbf{x}^{(h)}), \quad (2)$$

which is illustrated in Fig. 1(a), where  $\dim(\mathbf{x}^{(h+1)}) + \dim(\mathbf{z}^{(h)}) = \dim(\mathbf{x}^{(h)})$ , and the RG transformation  $R_h$  can be made invertible. At each level, the transformation  $R_h$  is a local bijective map,

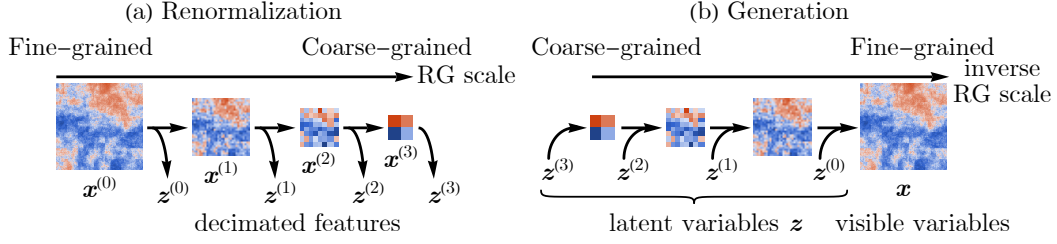


Figure 1: (a) The forward RG transformation splits out decimated features at different scales. (b) The inverse RG transformation generates the fine-grained image from latent variables.

which is constructed by stacking trainable bijective blocks. We will specify its details later. The split-out information  $\mathbf{z}^{(h)}$  can be viewed as latent variables arranged at different scales. Then the inverse RG transformation  $G_h \equiv R_h^{-1}$  simply generates the fine-grained image,

$$\mathbf{x}^{(h)} = R_h^{-1}(\mathbf{x}^{(h+1)}, \mathbf{z}^{(h)}) = G_h(\mathbf{x}^{(h+1)}, \mathbf{z}^{(h)}). \quad (3)$$

The highest-level image  $\mathbf{x}^{(h_L)} = G_{h_L}(\mathbf{z}^{(h_L)})$  can be considered as generated directly from latent variables  $\mathbf{z}^{(h_L)}$  without referring to any higher-level coarse-grained image, where  $h_L = \log_2 L - \log_2 m$ , for the original image of size  $L \times L$  with local transformations acting on kernel size  $m \times m$ . Therefore, given the latent variables  $\mathbf{z} = \{\mathbf{z}^{(h)}\}$  at all levels  $h$ , the original image can be restored by the following nested maps, as illustrated in Fig. 1(b),

$$\mathbf{x} \equiv \mathbf{x}^{(0)} = G_0(G_1(G_2(\dots, \mathbf{z}^{(2)}), \mathbf{z}^{(1)}), \mathbf{z}^{(0)}) \equiv G(\mathbf{z}), \quad (4)$$

where  $\mathbf{z} = \{\mathbf{z}^{(0)}, \dots, \mathbf{z}^{(h_L)}\}$ . RG-Flow is a flow-based generative model that uses the above composite bijective map  $G$  as the generator.

To model the RG transformation, we arrange the bijective blocks in a hierarchical network architecture. Fig. 2(a) shows the side view of the network, where each green or yellow block is a local bijective map. Following the notation of MERA networks, the green blocks are the *disentangled*, which reparametrize local variables to reduce their correlations, and the yellow blocks are the *decimators*, which separate the decimated features out as latent variables. The blue dots on the bottom are the visible variables  $\mathbf{x}$  from the data, and the red crosses are the latent variables  $\mathbf{z}$ . We omit color channels of the image in the illustration, since we keep the number of color channels unchanged through the transformation.

Fig. 2(b) shows the *top-down view* of a step of the RG transformation. The green/yellow blocks (disentangled/decimators) are interwoven on top of each other. The covering area of a disentangler or decimator is defined as the kernel size  $m \times m$  of the bijector. For example, in Fig. 2(b), the kernel size is  $4 \times 4$ . After the decimator, three fourth of the degrees of freedom are decimated into latent variables (red crosses in Fig. 2(a)), so the edge length of the image is halved.

As a mathematical description, for the single-step RG transformation  $R_h$ , in each block  $(p, q)$  labeled by  $p, q = 0, 1, \dots, \frac{L}{2^h m} - 1$ , the mapping from  $\mathbf{x}^{(h)}$  to  $(\mathbf{x}^{(h+1)}, \mathbf{z}^{(h)})$  is given by

$$\begin{aligned} \left\{ \mathbf{y}_{2^h(mq + \frac{m}{2} + a, mq + \frac{m}{2} + b)}^{(h)} \right\}_{(a,b) \in \square_m^1} &= R_h^{\text{dis}} \left( \left\{ \mathbf{x}_{2^h(mq + \frac{m}{2} + a, mq + \frac{m}{2} + b)}^{(h)} \right\}_{(a,b) \in \square_m^1} \right) \\ \left\{ \mathbf{x}_{2^h(mp + a, mq + b)}^{(h+1)} \right\}_{(a,b) \in \square_m^2}, \left\{ \mathbf{z}_{2^h(mp + a, mq + b)}^{(h)} \right\}_{(a,b) \in \square_m^1 / \square_m^2} &= R_h^{\text{dec}} \left( \left\{ \mathbf{y}_{2^h(mp + a, mq + b)}^{(h)} \right\}_{(a,b) \in \square_m^1} \right), \end{aligned} \quad (5)$$

where  $\square_m^k = \{(ka, kb) \mid a, b = 0, 1, \dots, \frac{m}{k} - 1\}$  denotes the set of pixels in a  $m \times m$  square with stride  $k$ , and  $\mathbf{y}$  is the intermediate result after the disentangler but not the decimator. The notation  $\mathbf{x}_{(i,j)}^{(h)}$  stands for the variable (a vector of all channels) at the pixel  $(i, j)$  and at the RG level  $h$  (similarly for  $\mathbf{y}$  and  $\mathbf{z}$ ). The disentanglers  $R_h^{\text{dis}}$  and decimators  $R_h^{\text{dec}}$  can be any bijective neural network. Practically, We use the coupling layer proposed in the Real NVP networks (Dinh et al.,

2017) to build them, with a detailed description in Appendix A. By specifying the RG transformation  $R_h = R_h^{\text{dec}} \circ R_h^{\text{dis}}$  above, the generator  $G_h \equiv R_h^{-1}$  is automatically specified as the inverse transformation.

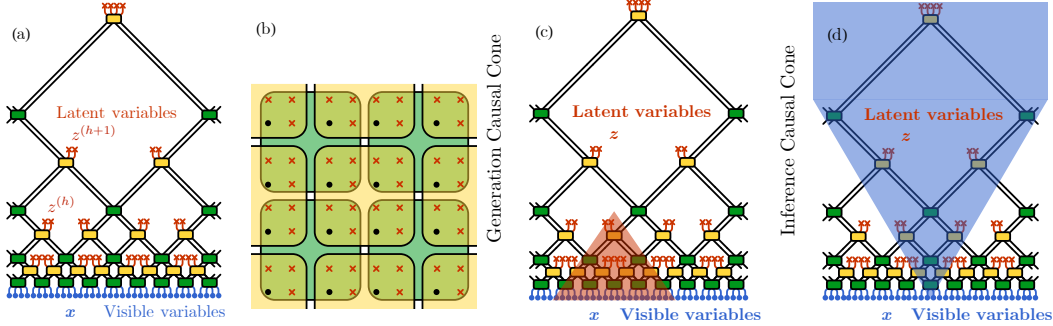


Figure 2: Subplot (a) shows the side view of the network. Green/Yellow blocks denote the disentangles/decimators, which are bijective maps. Subplot (b) shows the top-down view of the network. The red area in the subplot (c) illustrates the generation causal cone for a latent variable. The blue area in the subplot (d) illustrates the inference causal cone for a visible variable.

**Training objective.** After decomposing the statistics into multiple scales, we need to make the latent features decoupled. So we assume that the latent variables  $z$  are independent random variables, described by a factorized prior distribution

$$p_Z(z) = \prod_l p(z_l), \quad (6)$$

where  $l$  labels every element in  $z$ , including the RG level, the pixel position and the channel. This prior gives the network the incentive to minimize the mutual information between latent variables. This *minimal bulk mutual information* (minBMI) principle was previously proposed to be the information theoretic principle that defines the RG transformation (Li & Wang (2018); Hu et al. (2020)).

Starting from a set of independent latent variables  $z$ , the generator  $G$  should build up correlations locally at different scales, such that the multi-scale correlation structure can emerge in the resulting image  $x$  to model the correlated probability distribution of the data. To achieve this goal, we should maximize the log likelihood for  $x$  drawn from the data set. The loss function to minimize reads

$$\mathcal{L} = -\mathbb{E}_{x \sim p_{\text{data}}(x)} \log p_X(x) = -\mathbb{E}_{x \sim p_{\text{data}}(x)} \left( \log p_Z(R(x)) + \log \left| \frac{\partial R(x)}{\partial x} \right| \right), \quad (7)$$

where  $R(x) \equiv G^{-1}(x) = z$  denotes the RG transformation, which contains trainable parameters. By optimizing the parameters, the network learns the optimal RG transformation from the data.

**Receptive fields of latent variables.** Due to the nature of local transformations in our hierarchical network, we can define the generation causal cone for a latent variable to be the affected area when that latent variable is changed. This is illustrated as the red cone in Fig. 2(c).

To visualize the latent space representation, we define the *receptive field* for a latent variable  $z_l$  as

$$\text{RF}_l = \mathbb{E}_{z \sim p_Z(z)} \left| \frac{\partial G(z)}{\partial z_l} \right|_c, \quad (8)$$

where  $|\cdot|_c$  denotes the 1-norm on the color channel. The receptive field reflects the response of the generated image to an infinitesimal change of the latent variable  $z_l$ , averaged over  $p_Z(z)$ . Therefore, the receptive field of a latent variable is always contained in its generation causal cone. Higher-level latent variables have larger receptive fields than those of the lower-level ones. Especially, if the receptive fields of two latent variables do not overlap, which is often the case for lower-level latent variables, they automatically become disentangled in the representation.

**Image inpainting and error correction.** Another advantage of the network locality can be demonstrated in the inpainting task. Similar to the generation causal cone, we can define the *inference causal cone* shown as the blue cone in Fig. 2(d). If we perturb a pixel at the bottom of the blue cone, all the latent variables within the blue cone will be affected, whereas the latent variables outside the cone cannot be affected. An important property of the hyperbolic tree-like network is that the higher level contains exponentially fewer latent variables. Even though the inference causal cone is expanding as we go into higher levels, the number of latent variables dilutes exponentially as well, resulting in a constant number of latent variables covered by the inference causal cone on each level. Therefore, if a small local region on an image is corrupted, only  $O(\log L)$  latent variables need to be modified, where  $L$  is the edge length of the entire image. While for globally connected networks, all  $O(L^2)$  latent variables have to be varied.

**Sparse prior distribution.** We have chosen to hard-code the RG information principle by using a factorized prior distribution, i.e.  $p_Z(\mathbf{z}) = \prod_l p(z_l)$ . The common practice is to choose  $p(z_l)$  to be the standard Gaussian distribution, which is spherical symmetric. If we apply any rotation to  $\mathbf{z}$ , the distribution will remain the same. Therefore, we cannot avoid different features from being mixed under the arbitrary rotation.

To overcome this issue, we use an anisotropic sparse prior distribution for  $p_Z(\mathbf{z})$ . In our implementation, we choose the Laplacian distribution  $p(z_l) = \frac{1}{2b} \exp(-|z_l|/b)$ , which is sparser compared to Gaussian distribution and breaks the spherical symmetry of the latent space. In Appendix E, we show a two-dimensional pinwheel example to illustrate this intuition. This heuristic method will encourage the model to find more semantically meaningful representations by breaking the spherical symmetry.

## 4 EXPERIMENTS

**Synthetic multi-scale datasets.** To illustrate RG-Flow’s ability to disentangle representations at different scales and spatially separated representations, we propose two synthetic datasets with multi-scale features, named MSDS1 and MSDS2. Their samples are shown in Appendix B. In each image, there are 16 ovals with different colors and orientations. In MSDS1, all ovals in an image have almost the same color, while their orientations are randomly distributed. So the color is a global feature in MSDS1, and the orientation is a local feature. In MSDS2, on the contrary, the orientation is a global feature, and the color is a local one.

We implement RG-Flow as shown in Fig. 2. After training, we find that RG-Flow can easily capture the characteristics of those datasets. Namely, the ovals in each image from MSDS1 have almost the same color; and from MSDS2, the same orientation. Especially, in Fig. 3, we plot the effect of varying latent variables at different levels, together with their receptive fields. For MSDS1, if we vary a high-level latent variable, the color of the whole image will change, which shows that the network has captured the global feature of the dataset. And if we vary a low-level latent variable, the orientation of only the corresponding one oval will change. As the ovals are spatially separated, the low-level representation of different ovals is disentangled. Similarly, for MSDS2, if we vary a high-level latent variable, the orientations of all ovals will change. And if we vary a low-level latent variable, the color of only the corresponding one oval will change.

For comparison, we also trained Real NVP on our synthetic datasets. We find that Real NVP fails to learn the global and local characteristics of those datasets. Details can be found in Appendix B.

**Human face dataset.** Next, we apply RG-Flow to more complicated multi-scale datasets. Most of our experiments use the human face dataset CelebA (Liu et al., 2015), and we crop and scale the images to  $32 \times 32$  pixels. Details of the network and the training procedure can be found in Appendix A. Experiments on other datasets, such as CIFAR-10 (Krizhevsky et al.), and quantitative evaluations can also be found in Appendix G.

After training, the network learns to progressively generate finer-grained images, as shown in Fig. 4(a). The colors in the coarse-grained images are not necessarily the same as those at the same positions in the fine-grained images, because there is no constraint to prevent the RG transformation from mixing color channels.

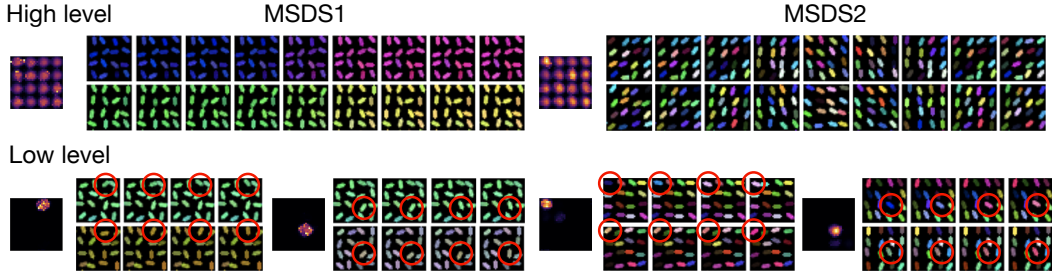


Figure 3: Multi-scale latent representations for MSDS1 and MSDS2.

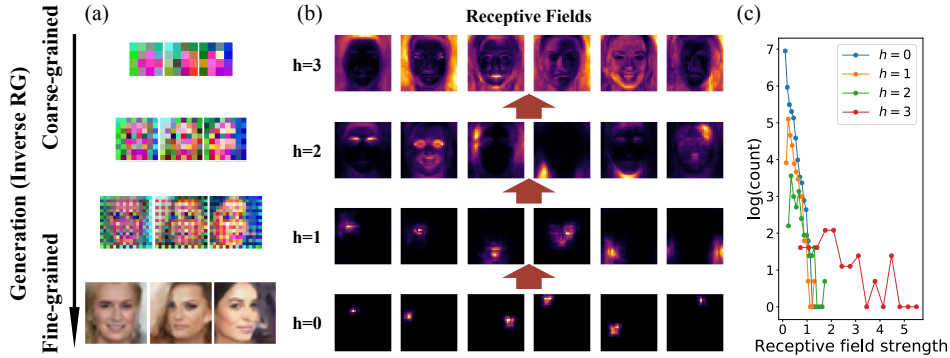


Figure 4: Subplot (a) shows the progressive generation of images during the inverse RG. Subplot (b) shows some receptive fields of latent variables from low level to high level. The strength of each receptive field is rescaled to one for better visualization. Subplot (c) shows the statistics of the receptive fields’ strength.

**Receptive fields.** To visualize the latent space representation, we calculate the receptive field for each latent variable, and list some of them in Fig. 4(b). We can see the receptive size is small for low-level variables and large for high-level ones, as indicated from the generation causal cone. In the lowest level ( $h = 0$ ), the receptive fields are merely small dots. In the second lowest level ( $h = 1$ ), small structures emerge, such as an eyebrow, an eye, a part of hair, etc. In the middle level ( $h = 2$ ), we can see eyebrows, eyes, forehead bang structure emerge. In the highest level ( $h = 3$ ), each receptive field grows to the whole image. We will investigate those explainable latent representations in the next section. For comparison, we show receptive fields of Real NVP in Appendix C. Even though Real NVP has multi-scale structure, since it is not locally constrained, semantic representations at different scales do not emerge.

**Learned features on different scales.** In this section, we show that some of these emergent structures correspond to explainable latent features. Flow-based generative model is the *maximal encoding procedure*, because the core of flow-based generative models is the bijective maps, and they preserves the dimensionality before and after the encoding. Usually, the images in the dataset live on a low dimensional manifold, and we do not need to use all the dimensions to encode such data. In Fig. 4(c) we show the statistics of the strength of receptive fields. We can see most of the latent variables have receptive fields with relatively small strength, meaning that if we change the value of those latent variables, the generated images will not be affected much. We focus on those latent variables with receptive field strength greater than one, which have visible effects on the generated images. We use  $h$  to label the RG level of latent variables, for example, the lowest-level latent variables have  $h = 0$ , whereas the highest-level latent variables have  $h = 4$ . In addition, we will focus on  $h = 1$  (low level),  $h = 2$  (mid level),  $h = 3$  (high level) latent variables. There are a few latent variables with  $h = 0$  that have visible effects, but their receptive fields are only small dots with no emergent structures.

For high-level latent representations, we found in total 30 latent variables that have visible effects, and six of them are identified with disentangled and explainable meanings. Those factors are gender, emotion, light angle, azimuth, hair color, and skin color. In Fig. 5(a), we plot the effect of varying those six high-level variables, together with their receptive fields. For the mid-level latent representations, we plot the four leading variables together with their receptive fields in Fig. 5(b), and they control eye, eyebrow, upper right bang, and collar respectively. For the low-level representations, some leading variables control an eyebrow and an eye as shown in Fig. 5(c). We see them achieve better disentangled representations when their receptive fields do not overlap.

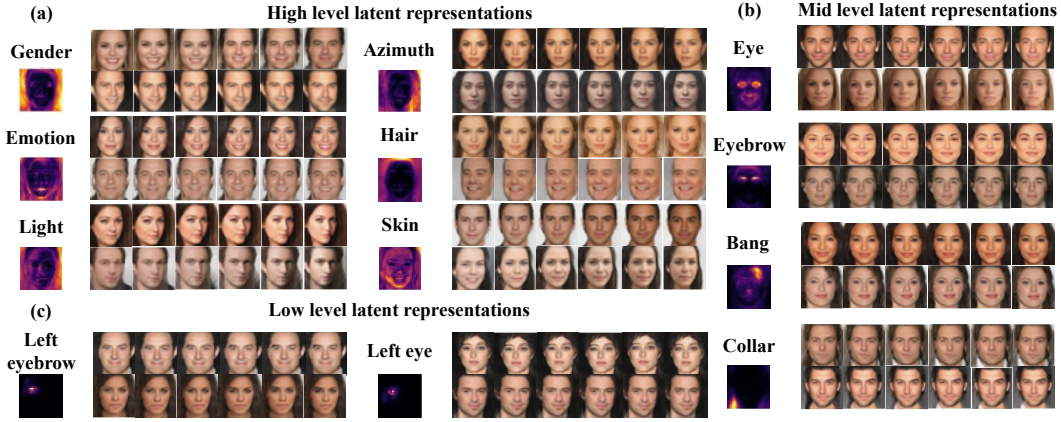


Figure 5: Semantic factors found on different levels.

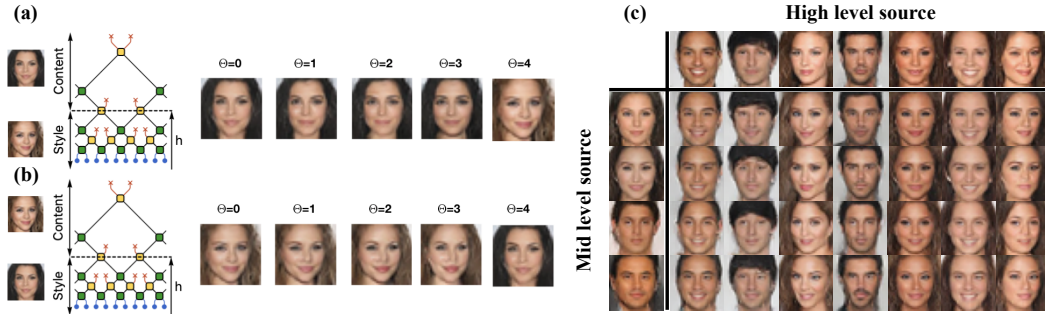


Figure 6: Image mixing in the hyperbolic tree-like latent space.

**Image mixing in scaling direction.** Given two images  $x_A$  and  $x_B$ , the conventional image mixing takes a linear combination between  $z_A = G^{-1}(x_A)$  and  $z_B = G^{-1}(x_B)$  by  $z = \lambda z_A + (1 - \lambda)z_B$  with  $\lambda \in [0, 1]$  and generates the mixed image from  $x = G(z)$ . In our model, latent variables  $z$  is coordinated by the pixel position  $(i, j)$  and the RG level  $h$ . The direct access of the latent variable  $z_{(i,j)}^{(h)}$  at each point enables us to mix the latent variables in a different manner, which may be dubbed as a “hyperbolic mixing”. We consider mixing the large-scale (high-level) features of  $x_A$  and the small-scale (low-level) features of  $x_B$  by combining their corresponding latent variables via

$$z^{(h)} = \begin{cases} z_A^{(h)}, & \text{for } h \geq \Theta, \\ z_B^{(h)}, & \text{for } h < \Theta, \end{cases} \quad (9)$$

where  $\Theta$  serves as a dividing line of the scales. As shown in Fig. 6(a), as we change  $\Theta$  from 0 to 3, more low-level information in the blonde-hair image is mixed with the high-level information of the black-hair image. Especially when  $h = 3$ , we see the mixed face have similar eyes, nose, eyebrows, and mouth as the blonde-hair image, while the high-level information, such as face orientation and hair color, is taken from the black-hair image. In addition, this mixing is not symmetric under the interchange of  $z_A$  and  $z_B$ , see Fig. 6(b) for comparison. This hyperbolic mixing achieves the similar



effect of StyleGAN (Karras et al., 2019; 2020) that we can take mid-level information from an image and mix it with the high-level information of another image. In Fig. 6(c), we show more examples of mixing faces.

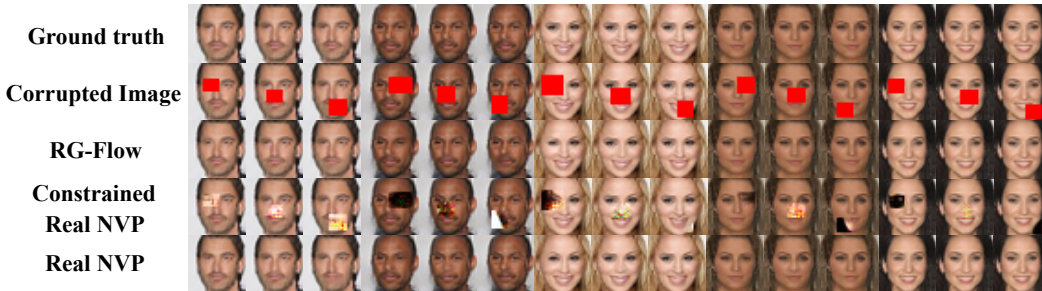


Figure 7: Inpainting locally corrupted images.

**Image inpainting and error correction.** The existence of the inference causal cone ensures that at most  $O(\log L)$  latent variables will be affected, if we have a small local corrupted region to be inpainted. In Fig. 7, we show that RG-Flow can faithfully recover the corrupted region (marked as red) only using latent variables locating inside the inference causal cone, which are around one third of all latent variables. For comparison, if we randomly pick the same number of latent variables to modify in Real NVP, it fails to inpaint as shown in Fig. 7 (Constrained Real NVP). To achieve the recovery of similar quality in Real NVP, as shown in Fig. 7 (Real NVP), all latent variables need to be modified, which are of  $O(L^2)$  order. See Appendix F for more details about the inpainting task and its quantitative evaluations.

## 5 DISCUSSION AND CONCLUSION

In this paper, we combined the ideas of renormalization group and sparse prior distribution to design RG-Flow, a probabilistic flow-based generative model. This versatile architecture can be incorporated with any bijective map to achieve an expressive flow-based generative model. We have shown that RG-Flow can separate information at different scales and encode them in latent variables living on a hyperbolic tree. To visualize the latent representations in RG-Flow, we defined the receptive fields for flow-based models in analogy to that in CNN. Taking CelebA dataset as our main example, we have shown that RG-Flow will not only find high-level representations, but also mid-level and low-level ones. The receptive fields serve as a visual guidance for us to find explainable representations. In contrast, the semantic representations of mid-level and low-level structures do not emerge in globally connected multi-scale flow models, such as Real NVP. We have also shown that the latent representations can be mixed at different scales, which achieves an effect similar to style mixing.

In our model, if the receptive fields of two latent representations do not overlap, they are naturally disentangled. For high-level representations, we propose to utilize a sparse prior to encourage disentanglement. We find that if the dataset only contains a few high-level factors, such as the 3D Chair dataset (Aubry et al., 2014) shown in Appendix G, it is hard to find explainable high-level disentangled representations, because of the redundant nature of the encoding in flow-based models. Incorporating information theoretic criteria to disentangle high-level representations in the redundant encoding procedure will be an interesting future direction.

## REFERENCES

- Samuel K Ainsworth, Nicholas J Foti, Adrian KC Lee, and Emily B Fox. oi-vae: Output interpretable vaes for nonlinear group factor analysis. In *International Conference on Machine Learning*, pp. 119–128, 2018.
- Tetsuya Akutagawa, Koji Hashimoto, and Takayuki Sumimoto. Deep learning and ads/qcd. *Physical Review D*, 102(2), Jul 2020. ISSN 2470-0029. doi: 10.1103/physrevd.102.026020.

- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics*, 6:483–495, 2018.
- Mathieu Aubry, Daniel Maturana, Alexei A. Efros, Bryan C. Russell, and Josef Sivic. Seeing 3d chairs: Exemplar part-based 2d-3d alignment using a large dataset of CAD models. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014*, pp. 3762–3769. IEEE Computer Society, 2014.
- Jens Behrmann, Will Grathwohl, Ricky T. Q. Chen, David Duvenaud, and Jörn-Henrik Jacobsen. Invertible residual networks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019*, volume 97 of *Proceedings of Machine Learning Research*, pp. 573–582. PMLR, 2019.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Cédric Bény and Tobias J Osborne. The renormalization group via statistical inference. *New Journal of Physics*, 17(8):083005, aug 2015. doi: 10.1088/1367-2630/17/8/083005.
- Urs Bergmann, Nikolay Jetchev, and Roland Vollgraf. Learning texture manifolds with the periodic spatial GAN. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 469–477. PMLR, 2017.
- Johann Brehmer and Kyle Cranmer. Flows for simultaneous manifold learning and density estimation. *CoRR*, abs/2003.13913, 2020.
- Cédric Bény. Deep learning and the renormalization group, 2013.
- Tian Qi Chen, Xuechen Li, Roger B. Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018*, pp. 2615–2625, 2018a.
- Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018*, pp. 6572–6583, 2018b.
- Tian Qi Chen, Jens Behrmann, David Duvenaud, and Jörn-Henrik Jacobsen. Residual flows for invertible generative modeling. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019*, pp. 9913–9923, 2019.
- Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016*, pp. 2172–2180, 2016.
- Brian Cheung, Jesse A. Livezey, Arjun K. Bansal, and Bruno A. Olshausen. Discovering hidden factors of variation in deep networks. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- Iris Cong, Soonwon Choi, and Mikhail D. Lukin. Quantum convolutional neural networks. *Nature Physics*, 15(12):1273–1278, Dec 2019. ISSN 1745-2481. doi: 10.1038/s41567-019-0648-8.
- Emily L. Denton, Soumith Chintala, Arthur Szlam, and Rob Fergus. Deep generative image models using a laplacian pyramid of adversarial networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015*, pp. 1486–1494, 2015.
- James J DiCarlo and David D Cox. Untangling invariant object recognition. *Trends in cognitive sciences*, 11(8):333–341, 2007.
- Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: non-linear independent components estimation. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.

- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real NVP. In *5th International Conference on Learning Representations, ICLR 2017, Conference Track Proceedings*. OpenReview.net, 2017.
- Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. In *5th International Conference on Learning Representations, ICLR 2017*. OpenReview.net, 2017.
- Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Alex Lamb, Martín Arjovsky, Olivier Mastropietro, and Aaron C. Courville. Adversarially learned inference. In *5th International Conference on Learning Representations, ICLR 2017*. OpenReview.net, 2017.
- Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019*, pp. 7509–7520, 2019.
- G. Evenbly and G. Vidal. Class of highly entangled many-body states that can be efficiently simulated. *Phys. Rev. Lett.*, 112:240502, Jun 2014. doi: 10.1103/PhysRevLett.112.240502.
- Michael E. Fisher. Renormalization group theory: Its basis and formulation in statistical physics. *Rev. Mod. Phys.*, 70:653–681, Apr 1998. doi: 10.1103/RevModPhys.70.653.
- Andrew Gambardella, Atilim Günes Baydin, and Philip H. S. Torr. Transflow learning: Repurposing flow models without retraining. *CoRR*, abs/1911.13270, 2019.
- Wen-Cong Gan and Fu-Wen Shu. Holography as deep learning. *International Journal of Modern Physics D*, 26(12):1743020, Oct 2017. ISSN 1793-6594. doi: 10.1142/s0218271817430209.
- L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2414–2423, 2016.
- Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Texture synthesis using convolutional neural networks. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015*, pp. 262–270, 2015.
- Koji Hashimoto. AdS/CFT correspondence as a deep boltzmann machine. *Phys. Rev. D*, 99:106017, May 2019. doi: 10.1103/PhysRevD.99.106017.
- Koji Hashimoto, Sotaro Sugishita, Akinori Tanaka, and Akio Tomiya. Deep learning and holographic qcd. *Phys. Rev. D*, 98:106014, Nov 2018. doi: 10.1103/PhysRevD.98.106014.
- Koji Hashimoto, Hong-Ye Hu, and Yi-Zhuang You. Neural ode and holographic qcd, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016*, pp. 770–778. IEEE Computer Society, 2016.
- Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *5th International Conference on Learning Representations, ICLR 2017*. OpenReview.net, 2017.
- Irina Higgins, David Amos, David Pfau, Sébastien Racanière, Loïc Matthey, Danilo J. Rezende, and Alexander Lerchner. Towards a definition of disentangled representations. *CoRR*, abs/1812.02230, 2018.
- Emiel Hoogeboom, Rianne van den Berg, and Max Welling. Emerging convolutions for generative normalizing flows. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2771–2780. PMLR, 2019.
- Hong-Ye Hu, Shuo-Hui Li, Lei Wang, and Yi-Zhuang You. Machine learning holographic mapping by neural network renormalization group. *Phys. Rev. Research*, 2:023369, Jun 2020.

- Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4):411–430, 2000.
- Nikolay Jetchev, Urs Bergmann, and Roland Vollgraf. Texture synthesis with spatial generative adversarial networks. *CoRR*, abs/1611.08207, 2016.
- Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision - ECCV 2016 - 14th European Conference, Proceedings, Part II*, volume 9906 of *Lecture Notes in Computer Science*, pp. 694–711. Springer, 2016.
- Leo P. Kadanoff. Scaling laws for ising models near  $T_c$ . *Physics Physique Fizika*, 2:263–272, Jun 1966.
- Mahdi Karami, Dale Schuurmans, Jascha Sohl-Dickstein, Laurent Dinh, and Daniel Duckworth. Invertible convolutional flow. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019*, pp. 5636–5646, 2019.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019*, pp. 4401–4410. Computer Vision Foundation / IEEE, 2019.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020*, pp. 8107–8116. IEEE, 2020.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2654–2663. PMLR, 2018.
- Diederik P. Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018*, pp. 10236–10245, 2018.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014*, 2014.
- Thomas N. Kipf, Elise van der Pol, and Max Welling. Contrastive learning of structured world models. In *8th International Conference on Learning Representations, ICLR 2020*. OpenReview.net, 2020.
- Maciej Koch-Janusz and Zohar Ringel. Mutual information, neural networks and the renormalization group. *Nature Physics*, 14(6):578–582, Jun 2018. ISSN 1745-2481. doi: 10.1038/s41567-018-0081-4.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- Y. LeCun. A theoretical framework for back-propagation. 1988.
- Yann LeCun, Bernhard E. Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne E. Hubbard, and Lawrence D. Jackel. Handwritten digit recognition with a back-propagation network. In *Advances in Neural Information Processing Systems 2*, pp. 396–404. Morgan Kaufmann, 1989.
- Ching Hua Lee and Xiao-Liang Qi. Exact holographic mapping in free fermion systems. *Phys. Rev. B*, 93:035112, Jan 2016. doi: 10.1103/PhysRevB.93.035112.
- Patrick M. Lenggenhager, Doruk Efe Gökmen, Zohar Ringel, Sebastian D. Huber, and Maciej Koch-Janusz. Optimal renormalization group transformation from information theory. *Phys. Rev. X*, 10: 011037, Feb 2020. doi: 10.1103/PhysRevX.10.011037.
- Shuo-Hui Li and Lei Wang. Neural network renormalization group. *Phys. Rev. Lett.*, 121:260601, Dec 2018.

- Henry W. Lin, Max Tegmark, and David Rolnick. Why does deep and cheap learning work so well? *Journal of Statistical Physics*, 168(6):1223–1247, Sep 2017. ISSN 1572-9613. doi: 10.1007/s10955-017-1836-5.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pp. 3730–3738. IEEE Computer Society, 2015. doi: 10.1109/ICCV.2015.425. URL <https://doi.org/10.1109/ICCV.2015.425>.
- Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019*, volume 97 of *Proceedings of Machine Learning Research*, pp. 4114–4124. PMLR, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019*. OpenReview.net, 2019.
- Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, and Ian J. Goodfellow. Adversarial autoencoders. *CoRR*, abs/1511.05644, 2015.
- Pankaj Mehta and David J. Schwab. An exact mapping between the variational renormalization group and deep learning, 2014.
- Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607, 1996.
- Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research*, 37(23):3311–3325, 1997.
- Dan Oprisa and Peter Toth. Criticality & deep learning ii: Momentum renormalisation group, 2017.
- Fernando Pastawski, Beni Yoshida, Daniel Harlow, and John Preskill. Holographic quantum error-correcting codes: toy models for the bulk/boundary correspondence. *Journal of High Energy Physics*, 2015(6):149, Jun 2015. ISSN 1029-8479. doi: 10.1007/JHEP06(2015)149.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019*, pp. 8024–8035, 2019.
- Xiao-Liang Qi. Exact holographic mapping and emergent space-time geometry. *arXiv: High Energy Physics - Theory*, 2013.
- Aditya Ramesh, Youngduck Choi, and Yann LeCun. A spectral regularizer for unsupervised disentanglement. *CoRR*, abs/1812.01161, 2018.
- Scott E. Reed, Aäron van den Oord, Nal Kalchbrenner, Sergio Gomez Colmenarejo, Ziyu Wang, Yutian Chen, Dan Belov, and Nando de Freitas. Parallel multiscale autoregressive density estimation. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2912–2921. PMLR, 2017.
- Danilo Jimenez Rezende, George Papamakarios, Sébastien Racanière, Michael S. Albergo, Gurtej Kanwar, Phiala E. Shanahan, and Kyle Cranmer. Normalizing flows on tori and spheres. *CoRR*, abs/2002.02428, 2020.
- Tim Salimans and Diederik P. Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016*, pp. 901, 2016.

- Robin Tibor Schirrmester, Yuxuan Zhou, Tonio Ball, and Dan Zhang. Understanding anomaly detection with deep invertible networks through hierarchies of distributions and features. *CoRR*, abs/2006.10848, 2020.
- H. Eugene Stanley. Scaling, universality, and renormalization: Three pillars of modern critical phenomena. *Rev. Mod. Phys.*, 71:S358–S366, Mar 1999. doi: 10.1103/RevModPhys.71.S358.
- Brian Swingle. Entanglement renormalization and holography. *Phys. Rev. D*, 86:065007, Sep 2012. doi: 10.1103/PhysRevD.86.065007.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Re-thinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 2818–2826. IEEE Computer Society, 2016. doi: 10.1109/CVPR.2016.308. URL <https://doi.org/10.1109/CVPR.2016.308>.
- Joshua B. Tenenbaum and William T. Freeman. Separating style and content with bilinear models. *Neural Comput.*, 12(6):1247–1283, 2000.
- Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, and Victor S. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pp. 1349–1357. JMLR.org, 2016.
- Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. *CoRR*, abs/2007.03898, 2020.
- G. Vidal. Class of quantum many-body states that can be efficiently simulated. *Phys. Rev. Lett.*, 101:110501, Sep 2008. doi: 10.1103/PhysRevLett.101.110501.
- Kenneth G. Wilson. Renormalization group and critical phenomena. i. renormalization group and the kadanoff scaling picture. *Phys. Rev. B*, 4:3174–3183, Nov 1971.
- Yi-Zhuang You, Xiao-Liang Qi, and Cenke Xu. Entanglement holographic mapping of many-body localized system by spectrum bifurcation renormalization group. *Phys. Rev. B*, 93:104205, Mar 2016. doi: 10.1103/PhysRevB.93.104205.
- Yi-Zhuang You, Zhao Yang, and Xiao-Liang Qi. Machine learning spatial geometry from entanglement features. *Phys. Rev. B*, 97:045153, Jan 2018. doi: 10.1103/PhysRevB.97.045153.
- Juexiao Zhang, Yubei Chen, Brian Cheung, and Bruno A. Olshausen. Word embedding visualization via dictionary learning. *CoRR*, abs/1910.03833, 2019.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017.

# Appendices for RG-Flow

## A DETAILS OF THE NETWORK AND THE TRAINING PROCEDURE

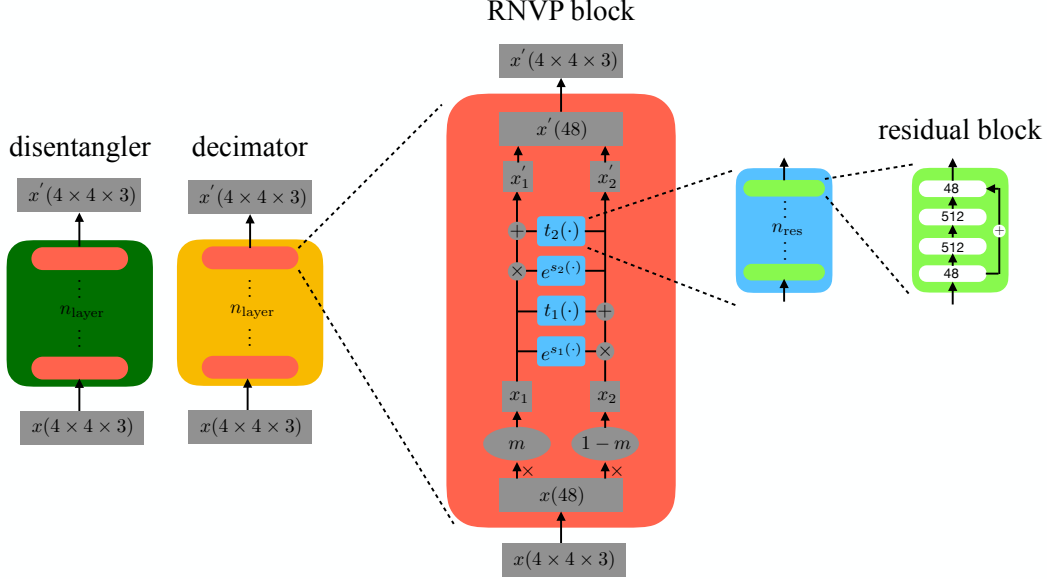


Figure 8: The details of the disentangler/decimator as a bijective map.

As shown in Fig. 2(a), each RG step includes a disentangler (green block) and a decimator (yellow block), which is parameterized by a bijective neural network. The blocks in every horizontal level share parameters, so we view them as a same block. In our reported experiments, disentanglers and decimators in different RG steps do not share parameters. However, we can also make them share parameters, which implements a scale-invariant course-graining process.

For each disentangler or decimator, we split the image into  $4 \times 4$  patches as shown in Fig. 2(b), stack them along the batch dimension, feed them into the network, and merge the output patches into a new image. After each RG step, the edge length of the image (the number of black lines above yellow blocks) is halved, except for the last (topmost) RG step that decimates all variables.

The choice of the bijective neural networks for the disentanglers and decimators can be versatile, and the performance of RG-Flow is strongly dependent on them. Since tuning the expressive power of those blocks is not the focus of our work, we choose the conventional Real NVP (Dinh et al., 2017) coupling layer to build disentanglers and decimators, which has great expressive power and is easy to invert. Fig. 8 illustrates our implementation. Each RNVP block is shown as a red block. It takes a  $4 \times 4$  image patch  $x$  as input, and split it into  $x_1$  and  $x_2$  using the checkerboard mask  $m$ . They are coupled to produce the output  $x'_2$ , using the formula

$$x'_2 = x_2 \odot \exp(s_1(x_1)) + t_1(x_1), \quad (10)$$

where  $\odot$  is the element-wise product.  $s(\cdot)$  and  $t(\cdot)$  are named the scale network and the translation network, which can be arbitrarily complicated to enhance the expressive power, as long as the numbers of input and output variables are the same. Then we use  $x'_2$  to alter  $x_1$  in the similar manner and outputs  $x'_1$ , and combine them to produce the output  $x'$  of the RNVP block.

We use residual networks with  $n_{\text{res}}$  residual blocks as  $s(\cdot)$  and  $t(\cdot)$ , shown as blue blocks in Fig. 8, and we choose  $n_{\text{res}} = 4$  in our implementation. Each residual block has 3 linear layers with size  $24 \rightarrow 512, 512 \rightarrow 512, 512 \rightarrow 24$ . Between linear layers, we insert swish activations (Chen et al., 2019), which is reported to give better results than ReLU and softplus, and its smoothness benefits our further analysis with higher order derivatives. We use Kaiming initialization (He et al., 2016) and weight normalization (Salimans & Kingma, 2016) on the linear layers.

The CelebA dataset contains rich information on different scales, including high-level information like gender and emotion, mid-level one like shapes of eyes and nose, and low-level one like details in hair and wrinkles. Because lower-level RG steps take larger images as input, we heuristically put more parameters in them. The numbers of RNVP blocks in all RG steps are  $n_{\text{layer}} = 8, 6, 4, 2$  respectively.

To preprocess the dataset, we use the aligned images from CelebA, crop a  $148 \times 148$  patch at the center of each image, downscale the patch to  $32 \times 32$  using bicubic downsampling, and randomly flip it horizontally.

We use AdamW optimizer (Loshchilov & Hutter, 2019) with conventional learning rate  $10^{-3}$  and weight decay rate  $5 \times 10^{-5}$ . To further stabilize training, we use gradient clipping with global norm 1. Between coupling layers, we use checkpointing (Paszke et al., 2019) to reduce memory usage. After using the checkpointing technique, we find the network’s memory consumption is greatly reduced. The maximal batch size can be set to 1024 on a single Nvidia Titan RTX given the current setup, which approximately has a million parameters. In our experiment, the batch size is conventionally set to 64. A training step takes about 1.2 seconds on an Nvidia Titan RTX.

The code for our implementation is at <https://hongyehu.github.io/RG-Flow/>

## B DETAILS OF THE SYNTHETIC MULTI-SCALE DATASETS

To illustrate RG-Flow’s ability to disentangle representations at different scales and spatially separated representations, we propose two synthetic toy datasets with multi-scale features, named MSDS1 and MSDS2, as shown in Fig. 9. Each dataset contains  $10^5$  images of  $32 \times 32$  pixels. In each image, there are 16 ovals with different colors and orientations, and their positions have small random variations to deform the  $4 \times 4$  grid. In MSDS1, all ovals in an image have almost the same color, while their orientations are randomly distributed. So the color is a global feature in MSDS1, and the orientation is a local feature. In MSDS2, on the contrary, the orientation is a global feature, and the color is a local one.

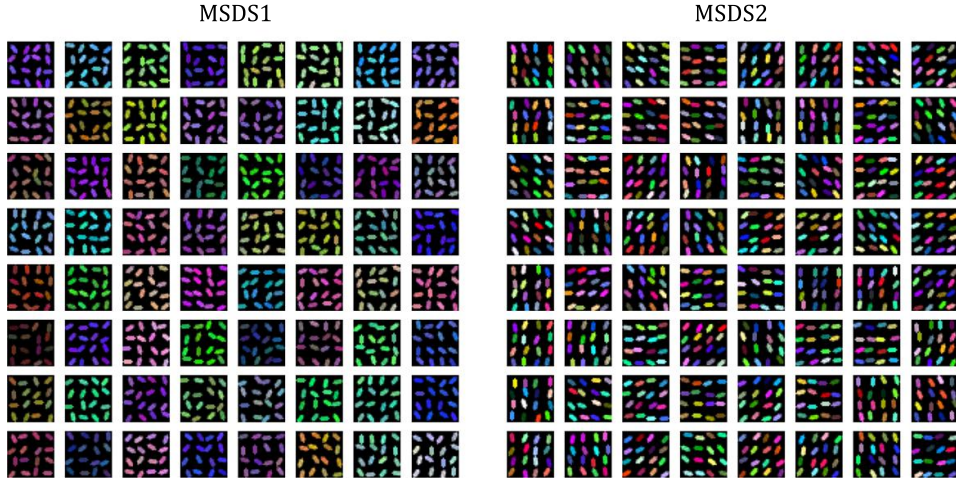


Figure 9: Samples from MSDS1 and MSDS2.

We trained RG-Flow on those datasets, with  $n_{\text{layer}} = 4$  and  $n_{\text{res}} = 4$ , and other hyperparameters described in Appendix. A. For comparison, we also trained Real NVP on those datasets, with approximately the same number of trainable parameters. Their generated images are shown in Fig. 10, where we can intuitively find that RG-Flow has learned the characteristics of the two datasets. Namely, in each image in MSDS1 the ovals should have almost the same color, and in MSDS2 the same orientation. In contrast, Real NVP fails to capture the global and local features of those datasets. The metrics of bits per dimension (BPD) and Fréchet Inception distance (FID) are listed in Table. 1. Note that FID may not reflect much semantic property for such synthetic datasets.



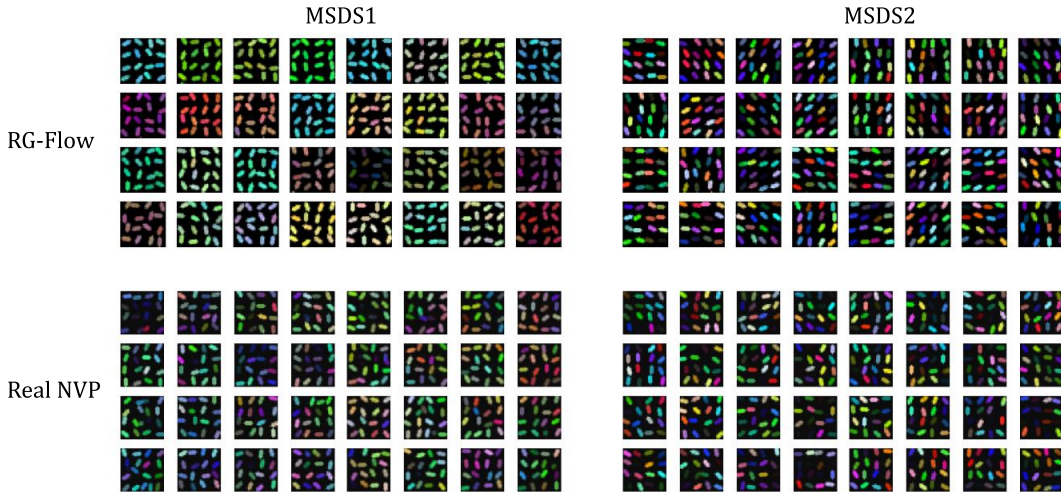


Figure 10: Samples from RG-Flow and Real NVP trained on MSDS1 and MSDS2.

	BPD ↓		FID ↓	
	MSDS1	MSDS2	MSDS1	MSDS2
RG-Flow	<b>0.906</b>	<b>1.01</b>	<b>1.40</b>	<b>2.61</b>
Real NVP	1.07	1.14	50.4	78.0

Table 1: BPD and FID from RG-Flow and Real NVP trained on MSDS1 and MSDS 2.

## C DISCUSSION ON REAL NVP

As a comparison to Fig. 4, we plot receptive fields of Real NVP in Fig. 11(a), together with the statistics of their strength in Fig. 11(b). Without local constraints on the bijective maps, there is no generation causal cone for Real NVP and other globally connected multi-scale models, and we cannot find semantic factors separated at different scales in general. In contrast, our RG-Flow can separate semantic factors at different scales, as shown in Fig. 4.

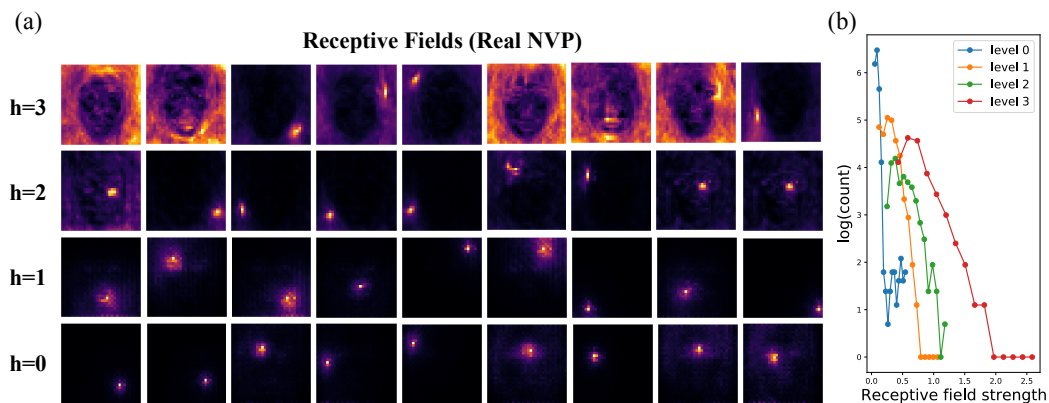


Figure 11: Subplot (a) shows the randomly picked receptive fields from the trained Real NVP on CelebA at different levels. Subplot (b) shows the statistics of the receptive fields' strength.

## D IMAGE GENERATION USING EFFECTIVE AND MIXED TEMPERATURES

Just like other generative models, we can define the effective temperature for RG-flow. Our prior distribution is

$$p_Z(\mathbf{z}) = \prod_l p(z_l), \quad (11)$$

where  $l = (i, j, h)$  labels every latent variable by its coordinate on the hyperbolic tree, especially  $h$  is the RG level. A model with effective temperature  $T$  ( $T > 0$ ) changes the prior distribution to  $p_{Z,T}(\mathbf{z})$ , with

$$p_{Z,T}(\mathbf{z}) \propto (p_Z(\mathbf{z}))^{1/T}. \quad (12)$$

For Laplacian prior, the effective temperature is implemented as

$$p_{Z,T}(\mathbf{z}) = \prod_l \frac{1}{2T} \exp\left(-\frac{|z_l|}{T}\right). \quad (13)$$

Moreover, we can define a mixed temperature model on the hyperbolic tree by

$$p_{Z,\text{mix}}(\mathbf{z}) \propto \prod_l (p(z_l))^{1/T_l}, \quad (14)$$

where  $T_l$  can be coordinate-dependent.

In our training procedure, we find the bijective maps on the lowest level take a longer time to converge. Therefore, when plotting the results in Fig. 5, we use the mixed temperature scheme with  $T_{h=0} = 0.2, T_{h=1} = T_{h=2} = T_{h=3} = 0.6$ , where  $T_h$  is the effective temperature on level  $h$ . Then we vary each latent variable from 0 to 1.5 if it is in high level or mid level, and from 0 to 6 if it is in low level.

## E A TOY MODEL FOR SPARSE PRIOR DISTRIBUTION

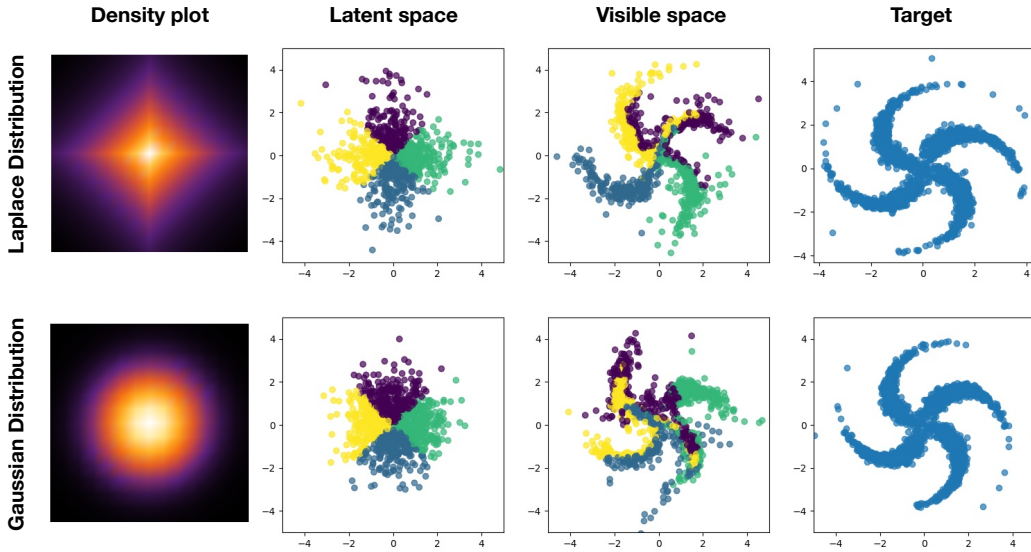


Figure 12: Two-dimensional pinwheel model.

In the first column of Fig. 12, the density plots of Laplacian distribution and Gaussian distribution are shown. As we can see, the Gaussian distribution has rotational symmetry, whereas the Laplacian distribution breaks  $SO(2)$  rotational symmetry to  $C_4$  symmetry. A flow-based generative model is trained to match the target distribution, which is a four-leg pinwheel as shown in the last column in Fig. 12. Given either Gaussian distribution or Laplacian distribution as the prior distribution, the model can learn the target distribution. In the second column, we sample 100 points and color them

by their living quadrant in the prior distribution. Then we map them to the visible space by the trained model, as shown in the third column. We see that four quadrants are approximately mapped to the four legs of the pinwheel if Laplacian prior is used. But for the Gaussian case, since it has rotational symmetry, the points in different quadrants are mixed more in the visible space, which makes it harder to interpret the mapping.

## F DETAILS OF INPAINTING EXPERIMENTS

For the inpainting experiments shown in Fig. 7, we randomly choose a corrupted region of  $10 \times 10$  pixels on the ground truth image, marked as the red patch in the second row of Fig. 7. We generate an image  $x_g$  from latent variables  $z_g$ , and use its corresponding region to fill in the corrupted region. Then we map the filled image  $x_f$  back to the latent variables  $z_f$ , and compute the log-likelihood. To recover the ground truth image, we optimize  $z_g$  to maximize the log-likelihood.

For RG-Flow, we only variate the latent variables living inside the inference causal cone, which is about 1200 out of 3072 latent variables. For the constrained Real NVP, we randomly pick the same amount of latent variables and allow them to variate, and we find it fails to inpaint the image in general. As a check, we find Real NVP can successfully inpaint the images if we optimize all latent variables, as shown in the last row of Fig. 7.

We use the conventional Adam optimizer to do the optimization. During the optimization procedure, we find the optimizer can be trapped in local minima. Therefore, for all experiments, we first randomly draw 200 initial samples of latent variables that are allowed to variate, then pick the one with the largest log-likelihood as the initialization.

To quantitatively assess the quality of inpainted images, we compute the peak signal-to-noise ratio (PSNR) of them against the ground truth images, and take the average over the 15 samples shown in Fig. 7. To further incorporate the semantic properties in the assessment, we can also use the Inception-v3 network (Szegedy et al., 2016). We first scale the images to  $299 \times 299$ , then feed them into the Inception-v3 network, extract the features from the “pool3” layer (as in FID score), and compute the PSNR. The results are listed in Table 2.

	PSNR $\uparrow$	Inception-PSNR $\uparrow$
RG-Flow	<b>38.7</b>	<b>24.9</b>
Constrained Real NVP	23.2	15.7
Real NVP	37.0	24.8

Table 2: PSNR and Inception-PSNR of inpainted images.

## G EXPERIMENTS ON OTHER DATASETS

For CIFAR-10 dataset, we use the same hyperparameters as described in Appendix. A. For 3D Chair dataset, we use  $n_{\text{layer}} = 8$  for all RG steps, because there is not so much low-level information as in CelebA and CIFAR-10. We also trained Real NVP on those datasets for comparison, with approximately the same number of trainable parameters. The metrics of bits per dimension (BPD) and Fréchet Inception distance (FID) are listed in Table. 3.

	BPD $\downarrow$			FID $\downarrow$		
	CelebA	CIFAR-10	3D Chair	CelebA	CIFAR-10	3D Chair
RG-Flow	3.47	<b>3.35</b>	<b>0.930</b>	31.3	<b>98.5</b>	<b>36.2</b>
Real NVP	<b>3.36</b>	3.61	0.933	<b>11.1</b>	126	73.8

Table 3: Bits per dimension (BPD) from RG-Flow and Real NVP trained on various datasets.

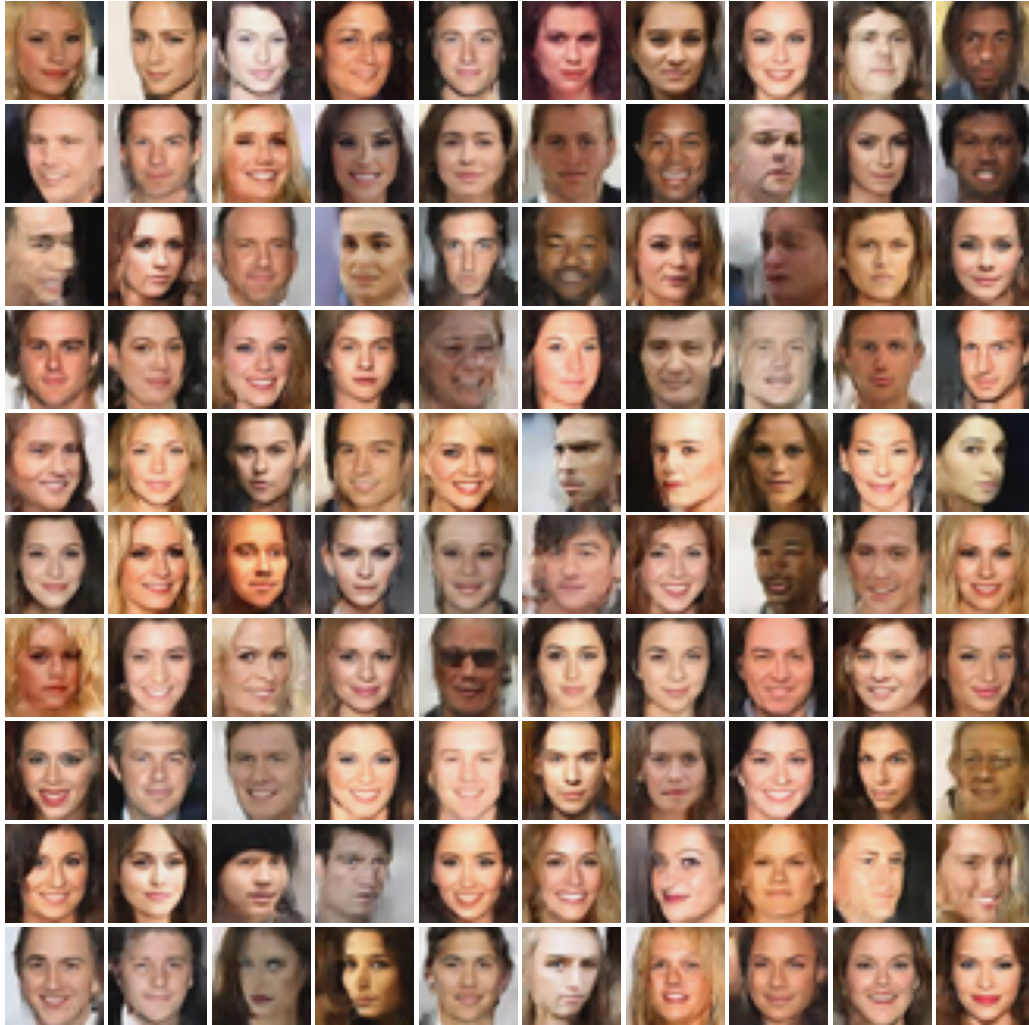


Figure 13: Samples from RG-Flow trained on CelebA dataset. We use  $T = 0.9$  when sampling.



Figure 14: Samples from RG-Flow trained on CIFAR-10 dataset.



Figure 15: Samples from RG-Flow trained on 3D Chair dataset.

## H RECEPTIVE FIELDS OF THE LATENT REPRESENTATIONS

More examples of randomly picked receptive fields are plotted in Fig. 16, Fig. 17, Fig. 18, and Fig. 19. For better visualization, we normalize all receptive fields' strength to one.

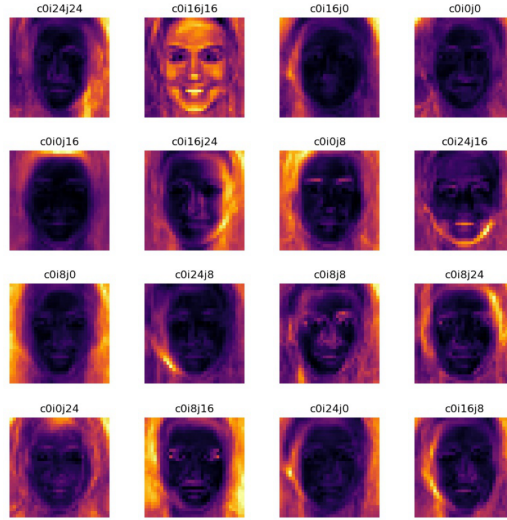


Figure 16: Receptive fields of high-level latent variables ( $h = 3$ ).

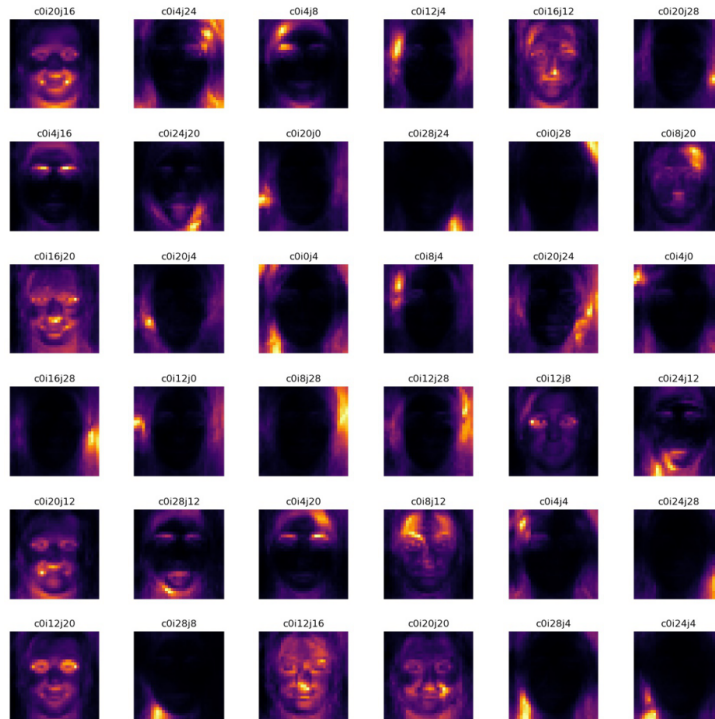
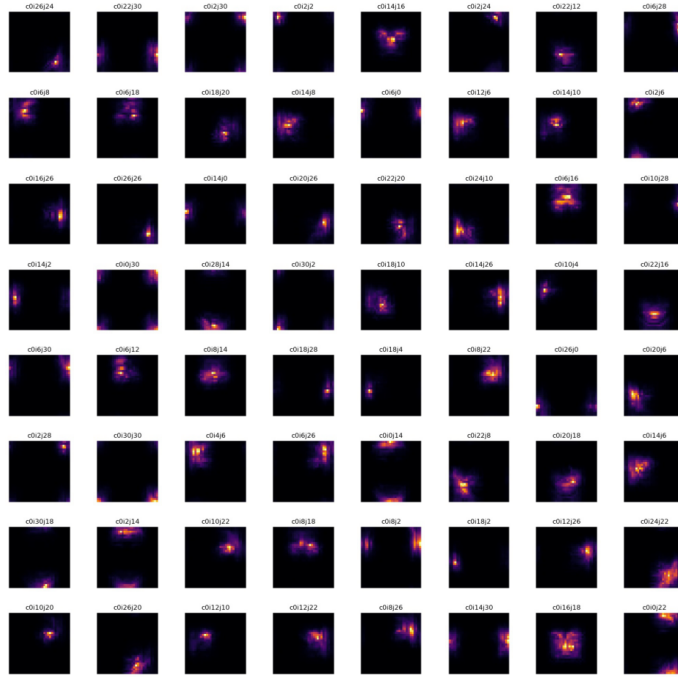
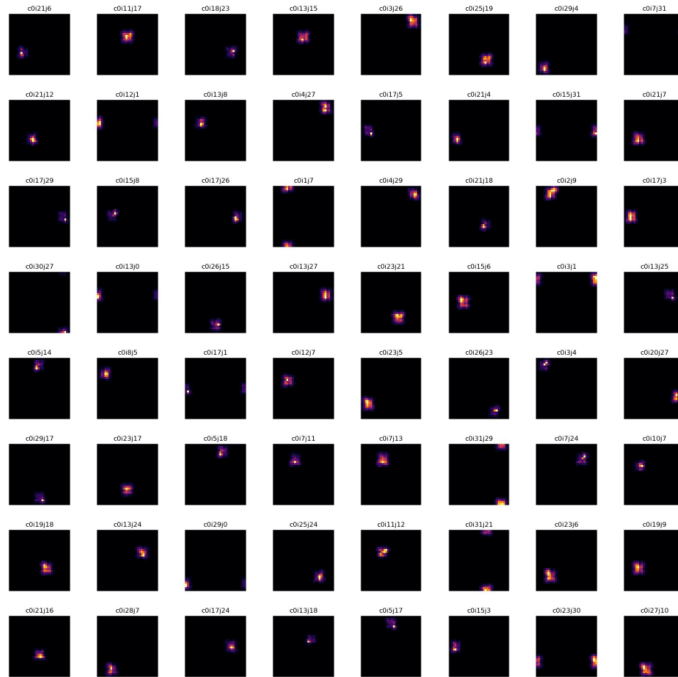


Figure 17: Receptive fields of mid-level latent variables ( $h = 2$ ).

Figure 18: Receptive fields of low-level latent variables ( $h = 1$ ).Figure 19: Receptive fields of the lowest-level latent variables ( $h = 0$ ).