

WHAT SECRETS DO YOUR MANIFOLDS HOLD? UNDERSTANDING THE LOCAL GEOMETRY OF GENERATIVE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep Generative Models are frequently used to learn continuous representations of complex data distributions using a finite number of samples. For any generative model, including pre-trained foundation models with GAN, Transformer or Diffusion architectures, generation performance can vary significantly based on which part of the learned data manifold is sampled. In this paper we study the post-training local geometry of the learned manifold and its relationship to generation outcomes for DDPM, Diffusion Transformer (DiT), an unconditional latent diffusion model and near state-of-the-art Stable Diffusion 1.4 text-to-image model. Building on the theory of continuous piecewise-linear (CPWL) generators, we characterize the local geometry in terms of three geometric descriptors - scaling (ψ), rank (ν), and complexity/un-smoothness (δ). We provide quantitative and qualitative evidence showing that for a given latent, the local descriptors are indicative of generation aesthetics, artifacts, diversity, and memorization. Finally we demonstrate that training a reward model using the local geometry of a pre-trained model, allows us to control the log-likelihood of a generated sample under the learned distribution and qualitative aspects of the generated image.

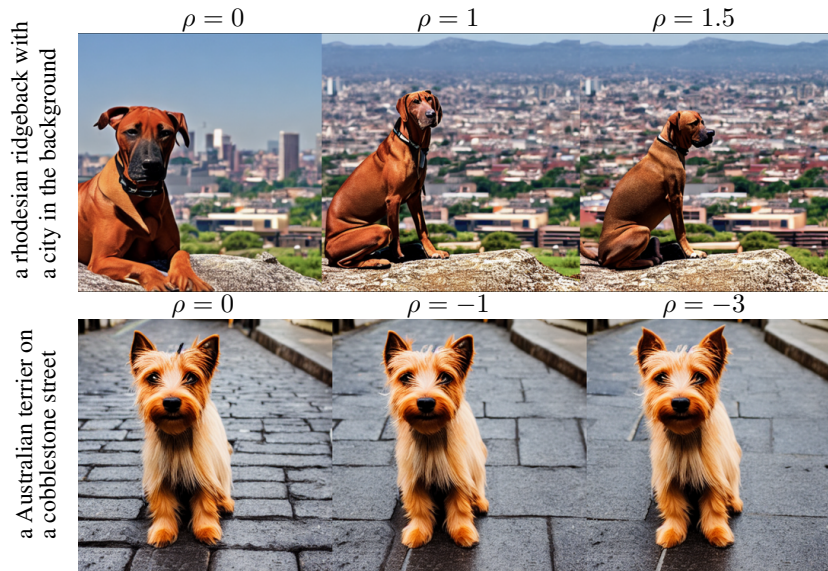


Figure 1: **Controlling visual complexity using geometry guidance.** We train a reward model on the geometric descriptor *local scaling* computed for the decoder of Stable Diffusion (Rombach et al., 2021). Positive (**top-row**) or negative (**bottom-row**) guidance ρ on this reward model allows decreasing (**top-row**) or increasing (**bottom-row**) the likelihood of the generated samples under the learned distribution of the Stable Diffusion decoder. As we decrease likelihood more background elements come into view and the focus on the subject decreases, vice-versa when increasing the likelihood.

1 INTRODUCTION

In recent years, deep generative models have emerged as a powerful tool in machine learning, capable of synthesizing realistic data across diverse domains (Karras et al., 2019; 2020; Rombach et al., 2021). However, the performance of such generative models may not be uniform across all downstream tasks, for example, recent studies have demonstrated that models like Stable Diffusion can exhibit biases in generation in terms of reduced generation fidelity or diversity, for certain demographic groups Zhao et al. (2018); Luccioni et al. (2023). Especially for 1) models trained with large heterogeneous data distributions, and 2) any pre-trained generative model without access to the training distribution, the aforementioned observations can become hard to interpret or reason. In such cases it is imperative to look at the model internals to obtain reasoning for model behavior. In this regard, we pose the following research question:

Research Question. *For any sample generated using a deep generative model, how is the local geometry of the model connected to downstream generation?*

The theory of continuous piecewise-linear generators Balestrieri et al. (2020) suggests that a large class of generative models can be considered continuous piecewise linear (CPWL) operators, implying that such generative models can be fully characterized in terms of their weights and architecture. We consider it the framework of choice to find answers to the aforementioned question and propose using three *local geometric descriptors* to quantify local characteristics of any pre-trained generative model:

- **Local rank** (ν), that characterizes the local dimensionality of the learned manifold.
- **Local scaling** (ψ), that characterizes the local change of volume by the generative model input output mapping.
- **Local complexity** (δ), that approximates the *un-smoothness* of the generative model in terms of second order changes in the input-output mapping.

Geometric descriptors such as local scaling, complexity or rank, have previously been used to measure function complexity of Deep Neural Networks (DNN) (Hanin & Rolnick, 2019) and DNN expressivity (Poole et al., 2016; Raghu et al., 2017), to evaluate the quality of representations learned with a self-supervised objective (Garrido et al., 2023), for interpretability and visualization of DNNs, (Humayun et al., 2023), to understand the learning dynamics in reinforcement learning (Cohan et al., 2022), to explain grokking, i.e., delayed generalization and robustness in classifiers (Humayun et al., 2024), maximum entropy or controlled sampling of GAN based generative models (Humayun et al., 2021; 2022b), and maximum likelihood inference in the latent space (Kuhnel et al., 2018). To the best of our knowledge, we are the first to use such geometric descriptors in the context of foundational scale text-to-image generative models.

Our contributions. In this paper, through rigorous experiments on large image generative models, we establish correlations between the local geometric descriptors and the aesthetic qualities, diversity, and degree of memorization of generated samples. We demonstrate how these manifest differently for different sub-populations of the generative distribution. We also show that the geometry of the data manifold is heavily influenced by the training data which enables applications in out-of-distribution detection and reward modeling to control the output distribution. Our empirical results lead to the following conclusions:

- **C1.** We present the first large-scale analysis of the local geometry of foundational text-to-image latent diffusion models and establish correlations between the local geometric descriptors and downstream aesthetic quality, diversity, and memorization (Sec 4.).
- **C2.** For small diffusion models and foundational image generative models, we show that the local geometry on the generative model manifold is distinct from the off manifold geometry, and can help distinguish the domain of a generative model (Sec. 3).
- **C3.** By training a auxiliary model on the local geometry of Stable Diffusion, we present a novel framework for reward guidance on a diffusion model to increase/decrease sampling diversity or control aesthetic qualities (Sec 5).

2 LOCAL DESCRIPTORS OF GENERATIVE MODEL MANIFOLDS

We start by introducing the geometric descriptors we will use in our study and provide intuition on what aspect of the generative model manifold geometry each of the descriptors quantify.

2.1 CONTINUOUS PIECEWISE-LINEAR GENERATIVE MODELS

Consider a generative network \mathcal{G} , which can be the decoder of a Variational Autoencoder (VAE) (Kingma & Welling, 2013), the generator of a Generative Adversarial Network (GAN) (Goodfellow et al., 2014), or an unrolled denoising diffusion implicit model (DDIM) (Song et al., 2020). Suppose, $\mathcal{G} : \mathbb{R}^E \rightarrow \mathbb{R}^D$ is a deep neural network with L layers, input space dimensionality E and output space dimensionality D . For any such generator, if the layers comprise affine operations such as convolutions, skip-connections, or max/avg-pooling, and the non-linearities are continuous piecewise-linear (CPWL) such as leaky-ReLU Xu (2015), ReLU, or periodic triangle, then the generator is a continuous piecewise-linear operator (Balestriero & Baraniuk, 2018a; Humayun et al., 2023). This implies that the $\mathcal{G} : \mathbb{R}^E \rightarrow \mathbb{R}^D$ mapping can be expressed in terms of a subdivision of the input space into linear regions Ω with each region ω from the input domain being mapped to the output via an affine operation. The continuous data manifold or image of the generator $\text{Im}(\mathcal{G})$ can be written as the union of sets:

$$\text{Im}(\mathcal{G}) = \bigcup_{\omega \in \Omega} \{\mathbf{A}_\omega z + \mathbf{b}_\omega \forall z \in \omega\}, \quad (1)$$

where, Ω is the partition of the latent space \mathbb{R}^E into continuous piecewise-linear regions, \mathbf{A}_ω and \mathbf{b}_ω are the slope and offset parameters of the affine mapping from latent space vectors $z \in \omega$ to the data manifold. For the class of continuous piecewise-linear (CPWL) neural network based generative models, Ω , \mathbf{A}_ω , and \mathbf{b}_ω are functions of the neurons/parameters of the network. For a generator with L layers, \mathbf{A}_ω and \mathbf{b}_ω can be expressed in closed-form in terms of the weights and the region-wise activation pattern of neurons for each layer. We refer the readers to Lemma 1 of (Humayun et al., 2023) for details.

To help build intuition, without loss of generality let's consider a CPWL toy generator that is trained on a handcrafted task where the target function $f : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ is a mapping between \mathbb{R}^2 and a mixture of five gaussian functions. Since the learned function is a continuous piecewise-affine spline operator, we use SplineCAM (Humayun et al., 2023) to analytically compute the function learned by the generator and visualize the learned manifold, as well as the input space piecewise-linear partition learned by the generator in Fig. 2 middle-left and left. Each convex region ω bounded by the black lines, is mapped to $\text{Im}(\mathcal{G})$ via per region parameters as described in Equation 1. The input-output mapping operation by the generator is affine region-wise, therefore any given input space region can be *scaled, rotated or translated* with a continuity constraint between regions, while going from the input to the output. For CPWL generators there are three characteristics of the learned manifold that can be studied: *i) the affine scaling induced per region, ii) the number of dimensions that are retained after scaling, i.e., local dimensionality of the learned manifold, and iii) the local smoothness of the CPWL partition.* We now introduce local descriptors that can be used to characterize these quantities.

2.1.1 LOCAL SCALING, ψ

We first introduce local scaling as a target descriptor to be used in our study that measures the local scaling performed on a region ω by a CPWL generator.

Definition 1. For a CPWL manifold produced by generator \mathcal{G} , the *local scaling* ψ_ω is constant within each region ω , and measures the log-scaling of the volume induced by the affine slope \mathbf{A}_ω for all latents $z \in \omega$. Local scaling for ω is expressed as

$$\psi_\omega = \log(\sqrt{\det(\mathbf{A}_\omega^T \mathbf{A}_\omega)}) = \sum_i^k \log(\sigma_i) \mathbb{1}_{\{\sigma_i \neq 0\}}, \quad (2)$$

where, $\{\sigma_i\}_{i=0}^k$, are the non-zero singular values of \mathbf{A}_ω .

Referring back to the example in Fig. 2, each region on the CPWL manifold (middle-left) and in the input space (left) is colored by ψ_ω , with darker shades indicating higher ψ_ω . Suppose \mathcal{G} has a uniform latent distribution, meaning every region ω has a uniform probability density in the latent space.

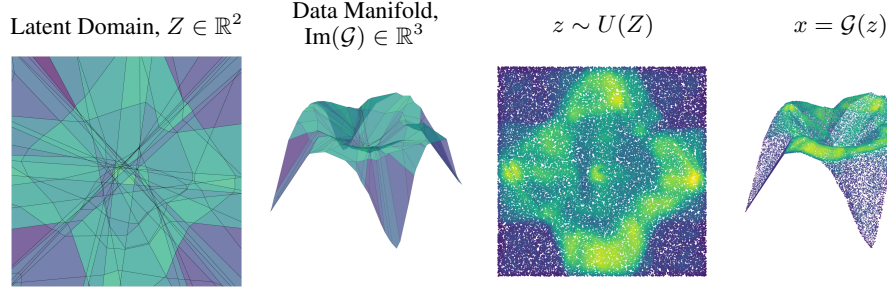


Figure 2: **The geometry of a continuous piecewise-linear toy generator.** For a CPWL generator $\mathcal{G} : \mathbb{R}^2 \rightarrow \mathbb{R}^3$, we provide analytically computed visualization of the input space partition, i.e., arrangement of linear regions (left) and learned CPWL manifold (middle-left). Each piece for this example, is colored by the piecewise-constant scaling induced by \mathcal{G} that is also analytically computed. Uniform samples from the latent domain (middle-right) and generated samples (right) are presented, colored by the estimated density at each sample using a gaussian kernel density estimator in \mathbb{R}^3 . We see that for any sample $z \in \omega$, the estimated density (\uparrow green) is inversely proportional to the scaling (\downarrow green) for region ω .

Under an injectivity assumption for any input space region ω and $S = \{\mathbf{A}_\omega \mathbf{z} + \mathbf{b}_\omega \forall \mathbf{z} \in \omega\}$, according to Theorem 1. in Humayun et al. (2022a), the output density on S , $p_S(\mathbf{x}) \propto \frac{1}{e^{\psi_\omega}}$. Therefore, local scaling ψ_ω is proportional to the negative log-likelihood of the generative model for any $z \in \omega$. We can validate this by using a kernel density estimator (KDE) to estimate the density of generated samples on the data manifold from a uniform latent distribution. In Fig. 2-right and middle-right, we denote the KDE estimated density per sample via colors, where higher density corresponds to brighter shades. Local scaling for any two regions $\omega \in \Omega$ and $\omega' \in \Omega$, can therefore be used to express the difference in generation uncertainty for two locations on the data manifold:

$$H_\omega - H_{\omega'} = \psi_\omega - \psi_{\omega'}, \quad (3)$$

where, $H_\omega, H_{\omega'}$ are the conditional entropy on the manifold for input space regions ω and ω' .

2.1.2 LOCAL RANK, ν .

The second descriptor we study is the rank of the region-wise slope matrix \mathbf{A}_ω , which represents the dimensionality of the manifold learned by a CPWL generator.

Definition 2. For a CPWL manifold produced by generator \mathcal{G} , *local rank* ν_ω is the exponent of the Shannon entropy of the spectral distribution of the per-region affine slope \mathbf{A}_ω and can be expressed as:

$$\nu_\omega = \exp \left(- \sum_i^k \alpha_i \log(\alpha_i) \right) \quad (4)$$

$$\text{where } \alpha_i = \frac{\sigma_i}{\sum_i^k \sigma_i} + \epsilon. \quad (5)$$

Here, $\{\sigma_i\}_{i=0}^{i=k}$ are non-zero singular values of \mathbf{A}_ω and $\epsilon = 10^{-30}$ is a constant. The local rank ν_ω can be shown to be equivalent to the dimensionality of the tangent space on the data manifold at \mathbf{z} .

2.1.3 LOCAL COMPLEXITY, δ

An important geometric notion to characterize any manifold locally is the local smoothness of the manifold. However, smoothness requires computing the hessian of the input-output mapping making it computationally intractable for large generative models. We therefore consider *local complexity* as a proxy for *sharpness* of the manifold locally for our study. Based on the notion of complexity for CPWL neural networks (Hanin & Rolnick, 2019), we can define local complexity of a CPWL generator as the following.

Definition 3. For a CPWL generator with input partition Ω , the *local complexity* δ_z for a P -dimensional neighborhood of radius r around latent vector z is

$$\delta_z = \sum_{\forall \omega \cap V_z \neq \emptyset} \mathbb{1}_\omega \quad (6)$$

$$\text{where } V_z = \{\mathbf{x} \in \mathbb{R}^E : \|\mathbf{B}(\mathbf{x} - z)\|_1 < r\}. \quad (7)$$

Here, \mathbf{B} is an orthonormal matrix of size $P \times E$ with $P \leq E$, $\|\cdot\|_1$ is the ℓ_1 norm operator and r is a radius parameter denoting the size of the locality to compute δ for. Here we consider a P dimensional neighborhood instead of the full dimensionality of the latent space to reduce computational complexity. The sum over regions $\omega \in V_z$ requires computing $\Omega \cap V_z$ which can be computationally intractable for high dimensions. A proxy for computing the partition for V_z with small r is counting the number of non-linearities within V_z , since for small r , the one can assume that the non-linearities do not fold inside V_z , therefore providing an upper bound on the number of regions according to Zaslavsky’s Theorem (Zaslavsky, 1975). To compute local complexity, we use the method introduced by Humayun et al. (2024) for general neural networks. We provide in appendix further implementation details and pseudocode.

2.2 EXTENDING BEYOND CONTINUOUS PIECEWISE-LINEAR GENERATORS

Networks with smooth activations. While the descriptors are defined for CPWL mappings, modern generative models employ a mixture of CPWL and non-CPWL operations. For networks with smooth activation functions or non-piecewise-linear non-linearities, our descriptors are first order Taylor approximations. We agree with the reviewer that there may be approximation errors incurred when we move from ReLU to smooth variants. However, Stable Diffusion already employs the GeLU activation function for which we perform the bulk of our experiments and find strong connections between the approximate local geometry and downstream generation. This is because smooth activation functions induce a soft VQ partitioning of the latent space compared to the hard VQ partitioning induced by a CPWL map Balestriero & Baraniuk (2018b). This suggest that much of the local linear structure we expect in CPWL maps are retained even if we employ smooth approximations of ReLU. Recent work has also empirically verified the local linearity for a large class of image based diffusion models Chen et al. (2024).

Computing jacobians for large networks. To avoid computing the singular values using the full input-output jacobian – which will be significantly expensive for large networks – when computing local scaling and rank we obtain singular values via randomized SVD Halko et al. (2011). First we obtain a random projection matrix with orthonormal rows \mathbf{W} with shape $k \times n$ such that $\mathbf{W}\mathbf{W}^T = \mathbb{I}_k$. Here n is the dimensionality of the outputs generated by the network. We therefore approximate local scaling as:

$$\psi_\omega^{(trunc)} = \sum_{i=1}^k \log(\sigma_i^{(trunc)}), \text{ where } \sigma_i^{(trunc)} \text{ are the non-zero singular values of } \mathbf{W}\mathbf{A}_\omega.$$

For any ω , if \mathbf{W} forms a basis for the range of \mathbf{A}_ω then $\sigma_i \approx \sigma_i^{(trunc)} \forall i = 1, 2, \dots, k$ [4]. Therefore $\mathbf{W}\mathbf{A}_\omega$ would provide us a low-rank approximation of \mathbf{A}_ω .

In our experiments we have tried two methods to obtain the projection matrix \mathbf{W} 1) by obtaining the eigenvectors for the covariance matrix for a set of 50K randomly generated samples. This was suggested in Halko et al. (2011). 2) by performing QR decomposition of a randomly initialized matrix. We see that the performance difference between methods 1) and 2) are negligible therefore consider the cheaper alternative 2) and consider a fixed pre-computed \mathbf{W} with $k=120$ for all \mathbf{A}_ω in our Stable Diffusion experiments.

3 CHARACTERIZING THE LOCAL GEOMETRY OF PRE-TRAINED MODELS VIA DESCRIPTORS

In this section, we explore the geometry of pre-trained generative models by characterizing the latent space to output manifold mapping in terms of the local geometric descriptors mentioned in the previous section. We are interested in the following questions: i) How does the on manifold local geometry vary from the off manifold local geometry? ii) How does the local geometry vary across the input domain?

3.1 ON AND OFF MANIFOLD GEOMETRY FOR DENOISING DIFFUSION PROBABILISTIC MODELS

Setup. To study the on and off manifold geometry of diffusion models, we train a denoising diffusion probabilistic model (DDPM) (Ho et al., 2020) on a toy dataset¹ to visualize how the local geometry varies for 1) different noise levels t , and 2) different training iterations. In Fig. 13-heatmaps we present the local complexity δ_x^t , local scaling ψ_x^t and local rank ν_x^t computed for different input space vectors x using the DDPM conditioned on noise levels t . Here T is the highest noise level in the forward diffusion process. We also present the difference between the expected descriptor values on and off the manifold, $\mathbb{E}_{\mathcal{M}}[\Phi] - \mathbb{E}_{\bar{\mathcal{M}}}[\Phi]$, $\forall \Phi \in \{\psi^t, \delta^t, \nu^t\}$ at different training iterations (right). We consider the set of input vectors within 0.05 units of the training data as on manifold \mathcal{M} and rest as off the manifold $\bar{\mathcal{M}}$.

Observations. The first observation is that with longer training, the maximum absolute difference between on and off manifold local geometry $\max_t \{|\mathbb{E}_{\mathcal{M}}[\Phi] - \mathbb{E}_{\bar{\mathcal{M}}}[\Phi]|\}$ increases. *Since with more training we see higher distinction between the on and off manifold geometry, this difference can be an indicator of learning in diffusion models.* We see that for well trained models, apart from $t > 0.17T$, ψ_x^t and ν_x^t decreases and δ_x^t increases with decreasing t , $\forall x \in \mathcal{M}$. This means, the likelihood on the manifold increases as noise levels are reduced, the smoothness decreases and the dimensionality of the manifold decreases as well. The quantity $\mathbb{E}_{\mathcal{M}}[\Phi] - \mathbb{E}_{\bar{\mathcal{M}}}[\Phi]$ is also minimized at $t \approx 0.17T$. This indicates that there can exist a noise level t conditioned on which diffusion model local scaling, rank and complexity have the highest distinction geometrically between on and off manifold vectors from the input space. *It can allow directly probing which parts of the input space are on the learned manifold* to possibly perform one step denoising or propose novel guidance schedules.

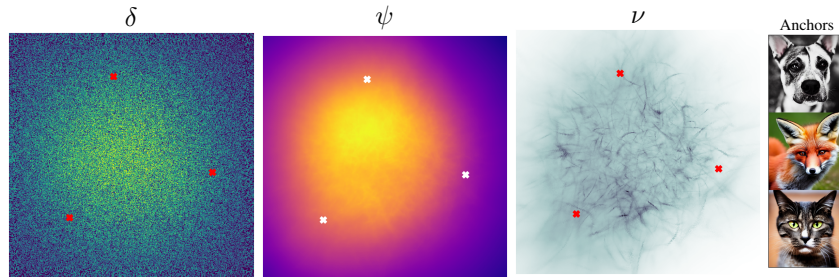


Figure 3: **Geometry of the Stable Diffusion latent space.** Geometric descriptors (left, middle-left, middle-right) visualized on a 2D latent space subspace, that passes through the latent representations of "a fox", "a cat" and "a dog" (right), denoted via markers on the 2D subspace descriptor. In Appendix, we provide denoised images for different high/low descriptor regions from the subspace. We see that in the convex hull of the three anchor latent vectors $\psi \uparrow$, $\nu \downarrow$ and $\delta \uparrow$. Moreover we see that in the convex hull, the local rank ν undergoes sharp changes which are not visible towards the edges of the domain.

3.2 THE LOCAL GEOMETRY OF LATENT DIFFUSION MODELS

In Sec. 3.1, we see that the local geometry in the input domain of a ddpm can be distinctive of its learned manifold. In this section we study the local geometry of the Stable Diffusion (SD) latent space, to explore whether there exists a relationship between the local geometry and the domain of the SD decoder.

Setup. While in Sec. 3.1 we could visualize the whole input domain of the diffusion model, for SD we can only visualize a subspace of the SD latent space. We use three prompts "a cat", "a dog" and "a fox" to generate three latent vectors using the SD diffusion model and consider a 2D slice in the latent space, going through the three denoised latents as our domain to visualize. Note that since this is a 2D subspace of the latent space, we can expect part of it to be in-domain for the SD decoder, whereas part of it would be out-of-domain. We provide implementation details in appendix.

Observations. We observe that 1) In the convex hull of the three denoised latents used as anchors for the 2D subspace being visualized, we have higher complexity, lower rank and higher local scaling. The decoded images from the convex hull may contain artifacts but are legible generations. 2) Local

¹<https://jumpingrivers.github.io/datasauRus>

rank does not smoothly vary across the latent space, especially with sharp changes in the local rank within the convex hull of the the anchor latent vectors. For the lowest rank regions in the convex hull, decoded images have good fidelity compared to latents with high uncertainty or complexity. 3) If we move away from the convex hull we see that generated images become more broken and contain heavy artifacts, indicating that such regions are out-of-domain for the SD decoder. However, we see that the local scaling is lower in these regions compared to the convex hull.

4 WHAT SECRETS DO YOUR MANIFOLDS HOLD?

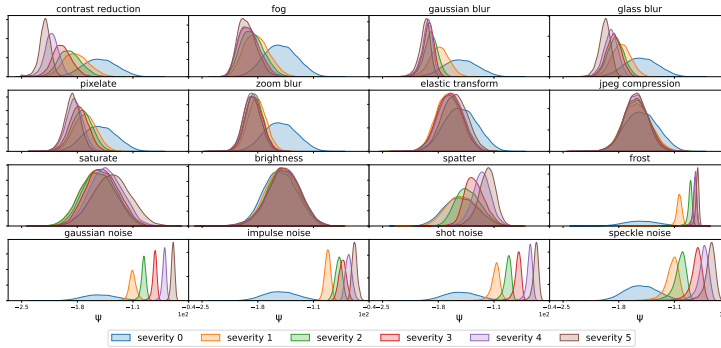


Figure 4: **Local scaling is sensitive to image corruptions.** We use 16 corruptions from (Hendrycks & Dietterich, 2019) to corrupt 10K imagenet images and compute the local geometry of the SD decoder. We see that SD local geometry is sensitive to corruptions, i.e., aesthetic changes to in-domain images (here Imagenet).

Correlations with visual complexity.

We selected 20K samples from Imagenet with resolution higher or equal to 512×512 , encode the samples using the SD encoder, and compute the local descriptors for the SD decoder. In Fig. 16 each column represents a local scaling level set, with the ψ for columns increasing from left to right. Recall that in Eq. 3 we show that increase in local scaling is equivalent to increase in uncertainty. In this figure, we can see that for lower uncertainty images we have more modal features in the images, i.e., the samples have less background elements and are focused on the subject corresponding to the Imagenet class.

For higher uncertainty images, we see that images have more outlier characteristics. For images with higher local rank in Fig. 18, we see that the backgrounds have higher frequency elements compared to lower rank images. For higher rank images, the dimensionality of the manifold is higher locally, therefore allowing more noise dimensions on the manifold. *These qualitative results provide evidence that the local geometry is indeed sensitive to natural variations of Imagenet images.*

Local geometry of corrupted images.

Fig. 6 illustrates the effect of applying 16 different image distortions (originally proposed in (Hendrycks & Dietterich, 2019)) to 10k ImageNet images. We consider Imagenet as in-domain for Stable Diffusion and encode them to the SD latent space to compute the geometric descriptors for the SD decoder. Samples are uniformly distributed over its classes. The plot shows the local scaling distribution at 6 increasing levels of severity $\in \{0, 1, 2, 3, 4, 5\}$, with zero corresponding to no corruptions applied. We observe that corruptions that are associated with reduction of spectral band, and/or reduction to the color range result in a reduction to the local scaling therefore the negative log-likelihood. We conjecture that this is due to the averaging effect of such distortions which move the corrupted images close to the mean of all images. Conversely, distortions known to be associated with the

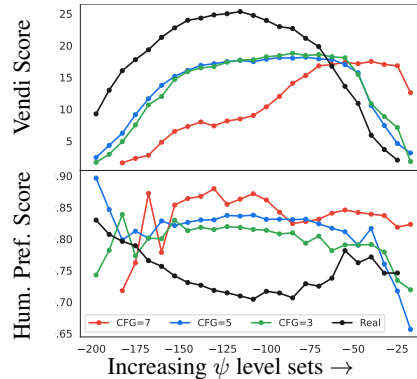


Figure 5: **Local Geometry level sets Imagenet prompts.** Vendi diversity scores and RAHF Liang et al. (2024) aesthetic scores computed for images with classifier free guidance (CFG) 7, 5 and 3. Diversity per level set increases and then decrease with increased local scaling. Aesthetic score slightly increases and then decreases as well with increased local scaling.

introduction of high-frequency artifacts are observed to produce an increase in local scaling therefore uncertainty moving the images away from the mean. *The results clearly indicate that the local geometry is sensitive to aesthetic changes to images introduced via most of the 16 corruptions.*

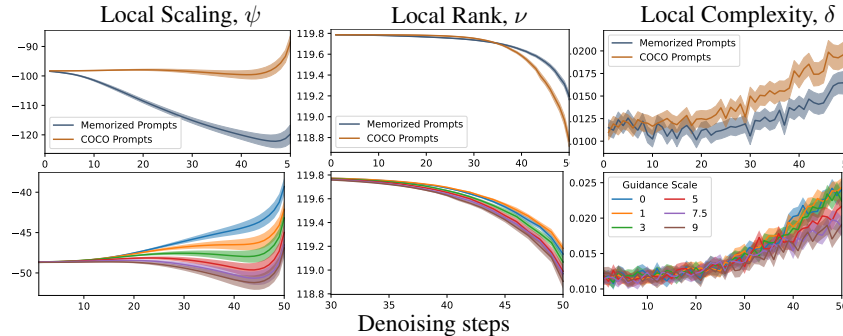


Figure 6: **Local geometry of denoising trajectories.** Geometric descriptors computed for the SD decoder unconditionally, during 50 stable diffusion denoising steps, for (top) 100 COCO and 100 memorized prompts (Wen et al., 2024) with guidance scale 7.5 and (bottom) 100 COCO prompts with varying guidance scales. For each prompt or guidance scale, we start from the same seeds. Shaded region represents 95% confidence interval. *We see that the local geometry trajectories are discriminative of memorization, as well as increased alignment when stronger classifier free guidance is used.*

Sensitivity to alignment and memorization.

Classifier-free-guidance is a method for increasing the alignment between the conditioning text-prompt and generated image Rombach et al. (2021). In Fig. 6-bottom, we see that for higher classifier free guidance scales during denoising, the avg. local scaling is lower, avg. local rank is lower and the avg. local complexity is higher. This indicates that for more conditionally aligned images obtained via classifier-free-guidance of the SD diffusion Unet, the decoder uncertainty is also lower especially during the final denoising steps. We also compute the local geometric descriptors for denoising trajectories conditioned on memorized prompts Wen et al. (2024) vs coco captions (Fig.6-top). We see that the mean local geometry is significantly different for denoising trajectories of 100 memorized prompts vs 100 random coco prompts. We also see that for higher guidance scales local rank ν is \downarrow and local complexity δ is \downarrow as well.

Downstream diversity and Human Preference Scores.

In Fig.8-left we present Vendi score aggregates and in middle and right, we present aesthetic and artifact score (higher is better) aggregates for real and generated Imagenet images with classifier free guidance scales of 7, 5 and 3, sorted in increasing local geometry level sets from left to right. Aesthetic and artifact scores are predicted by a human preference model Liang et al. (2024). We see that for increasing local scaling level sets, we have an increase in the diversity of images per bin. With very high local scaling, we get images from the highest uncertainty modes, i.e., the anti-modes, which result in a drop in the diversity of images. For real images, aesthetic and artifact scores get reduced for higher local scaling non-monotonically. This behavior is expected, as we have discussed before, higher local scaling images have higher uncertainty and higher visual complexity which may result in lower human preference predictions. For images from the highest local scaling bins however, we see an increase in the aesthetic and artifact scores. This is a unique phenomenon showing that for very uncertain images there could be a possible positive human preference. For local rank there is an increase in human preference scores for the high rank bins.

5 GUIDING GENERATION WITH GEOMETRY

In the previous sections, we have presented qualitative and quantitative evidence, establishing the connection between geometric descriptors and downstream generation. Among the three descriptors we find that local complexity has the highest sensitivity to aesthetic changes in images due to corruptions (section 4), and correlates with visual complexity (fig. 16). We also observe in section 4 higher local scaling level sets for samples generated using a classifier free guidance scale of 7.5, have

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

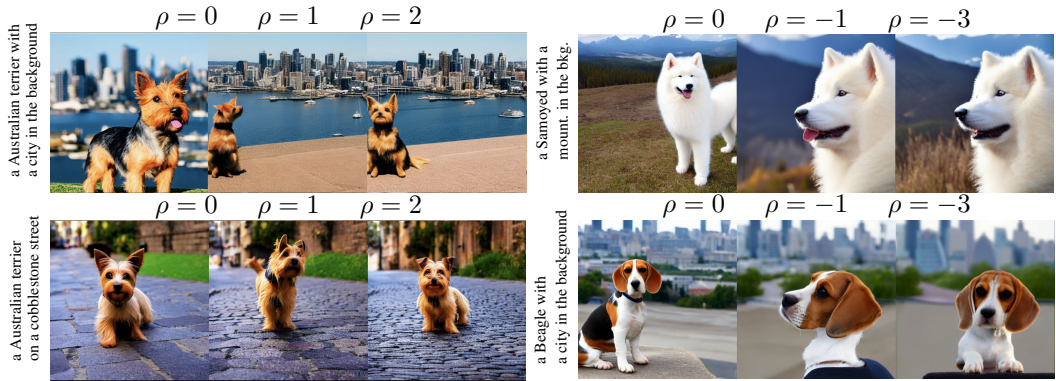


Figure 8: **Reward guidance on stable diffusion.** Local scaling reward guidance increases from left to right in each picture (left-panel) and decreases left to right in each picture (right-panel), with the first image showing no reward guidance. We observe maximizing the reward leads to sharper details, improved sharpness and contrast, and higher diversity in the images. Decreasing the local scaling leads to minimized uncertainty, resulting in a noticeable blurring effect and loss of details especially in the background of the image.

higher diversity while maintaining higher predicted human preference scores. Based on these results, we wish to explore whether local scaling can be used to guide generation.

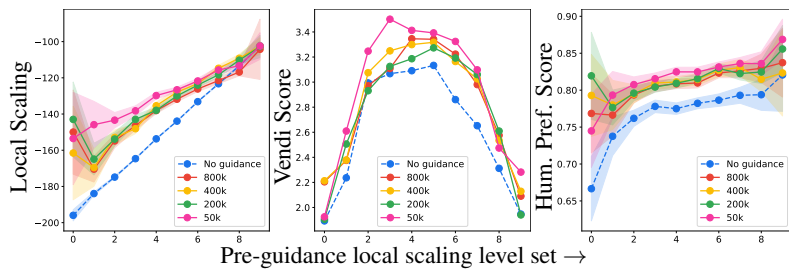


Figure 7: **Guidance via reward models trained with increasing number of training samples.** We observe that even with when only 50K samples are used in training, reward models can increase generation diversity, local scaling and human preference scores. Here we present the change of vendi score, average local scaling and average aesthetic score for pre-guidance local scaling level sets increasing from left to right on the x-axis.

that requires computing the input-output jacobian. To avoid computing the Hessian we train a reward model as a proxy and use the reward model gradients directly. Instead of training on continuous local scaling values in a regression task, we transform it into a local scaling level set classification task. We discretize the range of local scaling values into 5 bins and use the bin indices as training labels.

Data preparation. We obtain training data for the reward model by i) sampling N images from Imagenet and encoding them to the Stable Diffusion latent space ii) adding noise using the forward diffusion process up to randomly chosen noise levels iii) for each latent computing the local scaling descriptor.

To evaluate the performance of the reward model and dependency of the reward model on the number of training samples, we train multiple models for $N = 50K, 200K, 400K, 800K$. For evaluation, we generate 2560 samples using the dreambooth live subject prompt templates Ruiz et al. (2023), with Imgewoof Howard (2019) dogs as subjects. While Imgewoof dog classes are present in Imagenet therefore possibly in the training data, the dreambooth prompt templates contain a variety of settings that are not generally present in Imagenet, e.g., ‘a <subject> on top of pink fabric’.

Evaluation Setup. We first sample Stable diffusion without any reward guidance and with classifier-free guidance of 7.5 to obtain baseline samples. We partition the range of local scaling values obtained

Recently proposed instance-level universal guidance method (Bansal et al., 2023), can effectively influence the latents in the reverse process of a latent diffusion models to produce desired changes. Directly using local scaling to guide generation using such methods, require calculating the input-output Hessian since local scaling is a first-order measure

486 for the baseline samples into $n = 10$ bins fig. 7, where each bin contains images from a local scaling
487 level set. Following that we use the same seed and prompts as the baseline samples to generate
488 images using reward guidance to increase local scaling. For each bin or pre-guidance local scaling
489 level set, we compare between the baseline samples and corresponding reward guided generations
490 in the following three axes: i) change of local scaling ii) change of vendi (diversity) score iii) the
491 change of human preference score (RAHF Liang et al. (2024) aesthetic score).

492 **Results.** In fig. 7, we present the mean local scaling per bin with 95% confidence interval. Here the
493 blue line represents the mean pre-guidance local scaling values, increasing from left to right. In fig. 7,
494 we present vendi scores and average predicted human preference scores. For any bin, we present
495 results for the reward guidance scale that maximizes the local scaling. We see that even for a model
496 trained with 50K samples, we can have a considerable increase in the local scaling, diversity and
497 aesthetic score for most of the bins. Changes in local scaling, vendi and aesthetic scores are higher
498 for the lower pre-guidance local scaling level set bins compared to the higher pre-guidance local
499 scaling level set bins.

500 Our experiments reveal that maximizing local scaling in the manifold of a stable diffusion model
501 directly correlates with adding texture to the generated images. Moreover, this approach reduces
502 the likelihood on the manifold for single images. By optimizing the local scaling descriptor, the
503 generative model is guided towards producing more varied and textured outputs.

504 This approach is notable because traditional methods for diversity guidance generally function at the
505 distribution level. Our method, however, focuses on maximizing the inherent diversity as preserved
506 by the model within its learned manifold, effectively steering the generated images towards the
507 extremities of the distribution. This instance-level intervention allows for a more detailed and precise
508 enhancement of diversity, presenting a novel approach to guiding generative models.

509 As seen from Fig. 8 (left-panel) maximizing the reward results in added details in form of sharpening
510 the image, adding texture and contrast. We also observe that if we move towards minimizing the
511 reward, the images tend to loose fine-grained details as seen in Fig. 8 (right-panel). Please refer to the
512 supplementary material for more visual results.

514 6 CONCLUSION & FUTURE DIRECTIONS

515 In this paper, we present empirical evidence that the local geometric descriptors – local scaling
516 (ψ), local rank (ν) and local complexity (δ) - can effectively characterize the local geometry and
517 distinguish between downstream qualitative aspects of generated samples such as generation quality,
518 aesthetics, diversity, and memorization. Such descriptors only utilize the model’s architecture and
519 weights to characterize the behavior of generative models. We acknowledge two main limitations that
520 warrant further investigation. First, the geometry of the learned manifold is inherently influenced
521 by the training dynamics of the model. A deeper understanding of this relationship is needed to
522 fully leverage geometric analysis for models. Second, the computational complexity of our method,
523 particularly the calculation of the Jacobian matrix, may pose a practical challenge, especially for
524 large-scale models. Future work should explore more efficient algorithms or approximations to
525 address this limitation.

527 REFERENCES

- 528 Georgios Arvanitidis, Lars Kai Hansen, and Søren Hauberg. Latent space oddity: on the curvature of deep
529 generative models. *arXiv preprint arXiv:1710.11379*, 2017.
- 530 Randall Balestriero and Richard Baraniuk. A spline theory of deep learning. In *International Conference on*
531 *Machine Learning*, pp. 374–383. PMLR, 2018a.
- 532 Randall Balestriero and Richard G Baraniuk. From hard to soft: Understanding deep network nonlinearities via
533 vector quantization and statistical inference. *arXiv preprint arXiv:1810.09274*, 2018b.
- 534 Randall Balestriero, Sebastien Paris, and Richard Baraniuk. Max-affine spline insights into deep generative
535 networks. *arXiv preprint arXiv:2002.11912*, 2020.
- 536 Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and
537 Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on*
538 *Computer Vision and Pattern Recognition*, pp. 843–852, 2023.

540 Siyi Chen, Huijie Zhang, Minzhe Guo, Yifu Lu, Peng Wang, and Qing Qu. Exploring low-dimensional subspaces
541 in diffusion models for controllable image editing. *arXiv preprint arXiv:2409.02374*, 2024.

542

543 Setareh Cohan, Nam Hee Kim, David Rolnick, and Michiel van de Panne. Understanding the evolution of
544 linear regions in deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35:
545 10891–10903, 2022.

546 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural
547 information processing systems*, 34:8780–8794, 2021.

548 Dan Friedman and Adji Bousso Dieng. The vendi score: A diversity evaluation metric for machine learning.
549 *Transactions on Machine Learning Research*, 2023.

550 Quentin Garrido, Randall Balestriero, Laurent Najman, and Yann Lecun. Rankme: Assessing the downstream
551 performance of pretrained self-supervised representations by their rank. In *International Conference on
552 Machine Learning*, pp. 10929–10974. PMLR, 2023.

553

554 I. J Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.
555 Generative adversarial nets. In *Proceedings of the 27th International Conference on Neural Information
556 Processing Systems*, pp. 2672–2680. MIT Press, 2014.

557 Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness: Probabilistic
558 algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288, 2011.

559 Boris Hanin and David Rolnick. Complexity of linear regions in deep networks. *arXiv preprint arXiv:1901.09021*,
560 2019.

561 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions
562 and perturbations. In *International Conference on Learning Representations*, 2019. URL [https://
563 openreview.net/forum?id=HJz6tiCqYm](https://openreview.net/forum?id=HJz6tiCqYm).

564

565 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural
566 information processing systems*, 33:6840–6851, 2020.

567 Jeremy Howard. Imagewoof: a subset of 10 classes from imagenet that aren’t so easy to classify, March 2019.
568 URL <https://github.com/fastai/imagenette#imagewoof>.

569

570 Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. Magnet: Uniform sampling from deep
571 generative network manifolds without retraining. In *International Conference on Learning Representations*,
572 2021.

573 Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. MaGNET: Uniform sampling from deep
574 generative network manifolds without retraining. In *ICLR*, 2022a. URL [https://openreview.net/
575 forum?id=r5qumLiYwf9](https://openreview.net/forum?id=r5qumLiYwf9).

576 Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. Polarity sampling: Quality and diversity
577 control of pre-trained generative networks via singular values. In *CVPR*, pp. 10641–10650, 2022b.

578 Ahmed Imtiaz Humayun, Randall Balestriero, Guha Balakrishnan, and Richard G Baraniuk. Splinecam: Exact
579 visualization and characterization of deep network geometry and decision boundaries. In *Proceedings of the
580 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3789–3798, 2023.

581 Ahmed Imtiaz Humayun, Randall Balestriero, and Richard Baraniuk. Deep networks always grok and here is
582 why. *arXiv preprint arXiv:2402.15555*, 2024.

583

584 Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization in diffusion models
585 arises from geometry-adaptive harmonic representation. *arXiv preprint arXiv:2310.02557*, 2023.

586 Hamidreza Kamkari, Brendan Leigh Ross, Jesse C Cresswell, Anthony L Caterini, Rahul G Krishnan, and
587 Gabriel Loaiza-Ganem. A geometric explanation of the likelihood ood detection paradox. *arXiv preprint
588 arXiv:2403.18910*, 2024a.

589 Hamidreza Kamkari, Brendan Leigh Ross, Rasa Hosseinzadeh, Jesse C Cresswell, and Gabriel Loaiza-Ganem.
590 A geometric view of data complexity: Efficient local intrinsic dimension estimation with diffusion models.
591 *arXiv preprint arXiv:2406.03537*, 2024b.

592

593 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial
networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
4401–4410, 2019.

594 Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and
595 improving the image quality of stylegan. In *Proc. CVPR*, pp. 8110–8119, 2020.

596
597 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

598 Line Kuhnle, Tom Fletcher, Sarang Joshi, and Stefan Sommer. Latent space non-linear statistics. *arXiv preprint*
599 *arXiv:1805.07632*, 2018.

600 Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun, Jordi Pont-
601 Tuset, Sarah Young, Feng Yang, et al. Rich human feedback for text-to-image generation. In *Proceedings of*
602 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19401–19411, 2024.

603 Alexandra Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. Stable bias: Analyzing
604 societal representations in diffusion models. *arXiv preprint arXiv:2303.11408*, 2023.

605 Hila Manor and Tomer Michaeli. On the posterior distribution in denoising: Application to uncertainty
606 quantification. *arXiv preprint arXiv:2309.13598*, 2023.

607 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF*
608 *International Conference on Computer Vision*, pp. 4195–4205, 2023.

609 Ben Poole, Subhaneil Lahiri, Maithreyi Raghunathan, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential
610 expressivity in deep neural networks through transient chaos. In *Advances In Neural Information Processing*
611 *Systems*, pp. 3360–3368, 2016.

612 Maithra Raghunathan, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power
613 of deep neural networks. In *ICML*, pp. 2847–2854, 2017.

614 Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, and Yoshua Bengio. Contractive auto-encoders: Ex-
615 plicit invariance during feature extraction. In *Proceedings of the 28th international conference on international*
616 *conference on machine learning*, pp. 833–840, 2011.

617 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image
618 synthesis with latent diffusion models. 2022 IEEE. In *CVF Conference on Computer Vision and Pattern*
619 *Recognition (CVPR)*, pp. 10674–10685, 2021.

620 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth:
621 Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF*
622 *Conference on Computer Vision and Pattern Recognition*, pp. 22500–22510, 2023.

623 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint*
624 *arXiv:2010.02502*, 2020.

625 Yuxin Wen, Yuchen Liu, Chen Chen, and Lingjuan Lyu. Detecting, explaining, and mitigating memorization
626 in diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024. URL
627 <https://openreview.net/forum?id=84n3UwkH7b>.

628 Bing Xu. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*,
629 2015.

630 Thomas Zaslavsky. *Facing up to arrangements: Face-count formulas for partitions of space by hyperplanes:*
631 *Face-count formulas for partitions of space by hyperplanes*, volume 154. American Mathematical Soc., 1975.

632 Shengjia Zhao, Hongyu Ren, Arianna Yuan, Jiaming Song, Noah Goodman, and Stefano Ermon. Bias and
633 generalization in deep generative models: An empirical study. *arXiv preprint arXiv:1811.03259*, 2018.

634
635
636
637
638
639
640
641
642
643
644
645
646
647