

Image Generation with a Sphere Encoder

Kaiyu Yue^{†12} Menglin Jia^{†1} Ji Hou¹ Tom Goldstein²

Abstract

We introduce the Sphere Encoder, an efficient generative framework capable of producing images in a single forward pass and competing with many-step diffusion models using fewer than five steps. Our approach works by learning an encoder that maps natural images uniformly onto a spherical latent space, and a decoder that maps random latent vectors back to the image space. Trained solely through image reconstruction losses, the model generates an image by simply decoding a random point on the sphere. Our architecture naturally supports conditional generation, and looping the encoder/decoder a few times can further enhance image quality. Across several datasets, the sphere encoder approach yields performance competitive with state of the art diffusions, but with a small fraction of the inference cost. Project page is available at sphere-encoder.github.io.

1. Introduction

Most generative image models rely on either diffusion (Ho et al., 2020; Lipman et al., 2022) or autoregressive next-token prediction (Tian et al., 2024). With either paradigm, image generation is extremely slow and costly, requiring many forward passes to produce a single image.

We propose an alternative paradigm that is capable of generating sharp images with as little as one forward pass. Our approach, which we call a *sphere encoder*, works by training two complementary models: an encoder model that maps the distribution of natural images uniformly onto the sphere, and a decoder that maps points on the sphere back to natural images (Figure 2). The term aligns with the autoencoder convention, reflecting its encoder-decoder architecture. At test time, an image is generated quickly by sampling a random point on the sphere and passing it through the decoder.

Although the sphere encoder does not employ diffusion processes explicitly, it supports several key capabilities commonly associated with its diffusion-based cousins (Dhariwal & Nichol, 2021; Rombach et al., 2022; Esser et al., 2024).

[†]Equal contribution. ¹Meta ²University of Maryland.

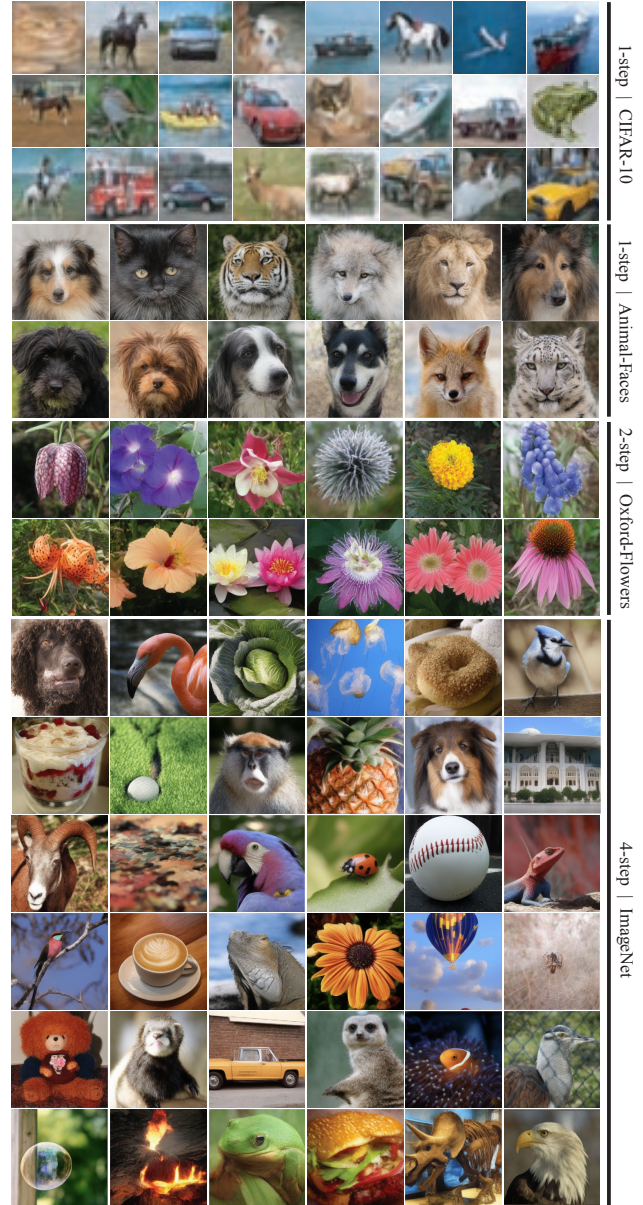


Figure 1. Selected images generated by the Sphere Encoder in one-step for CIFAR-10 (32×32) and Animal-Faces, two-steps for Oxford-Flowers, and four-steps for ImageNet (256×256).

These include conditional generation using AdaLN (Perez et al., 2018; Peebles & Xie, 2023), classifier-free guidance (CFG) (Ho & Salimans, 2022), and few-step iteration to enhance sample quality (Goodfellow et al., 2014; Kingma

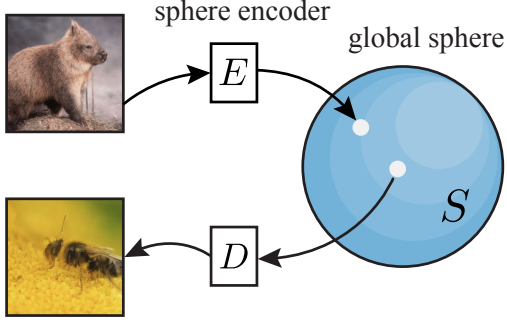


Figure 2. A sphere encoder E maps the natural image distribution uniformly onto a global sphere S . The decoder D then generates a realistic image by decoding a random point on the sphere.

& Dhariwal, 2018; Song et al., 2023). Experiments demonstrate that our approach achieves competitive one-step generation, and state-of-the-art performance in few-step regimes (e.g., fewer than 5 steps) on a range of datasets.

Motivation and Relation to Autoencoders

Autoencoders (LeCun, 1987; Bourlard & Kamp, 1988; Hinton & Zemel, 1993) have been widely used in representation learning and generative modeling. A lower-dimensional latent bottleneck between the encoder and decoder forces the model to learn an undercomplete representation of the input (Goodfellow et al., 2016).

To regularize the latent space, variational autoencoders (VAEs) (Kingma & Welling, 2013; 2019; Tolstikhin et al., 2018; Davidson et al., 2018; Ke & Xue, 2025) minimize the divergence between the latent distribution and a (typically) Gaussian prior. Unfortunately, in the standard VAE formulation, the divergence loss and image reconstruction loss are at odds with one another; zero divergence loss cannot be achieved simultaneously with perfect image reconstruction. As a result, the learned posterior fails to strongly match the prior – an issue known as the posterior hole problem (Makhzani et al., 2015; Rezende & Viola, 2018; Tomczak & Welling, 2018; Dai & Wipf, 2019; Ghosh et al., 2020; Aneja et al., 2021). Direct samples from the Gaussian prior fail to yield valid images. Realistic images are currently possible only by decoding samples from the posterior (i.e., adding noise to latents derived from real images), as illustrated in Figure 3. Our approach does not suffer from this problem.

Like a classical VAE, our approach relies on an autoencoder. Unlike the VAE, which tries to force the latent vectors into a Gaussian distribution, we instead force latents to be uniformly distributed on a sphere. Due to the bounded and rotationally symmetric nature of the sphere, this can be achieved simply by forcing embeddings of natural images away from one another, causing them to spread throughout the sphere. Moreover, this objective is not in contradiction with the image reconstruction objective; we can achieve

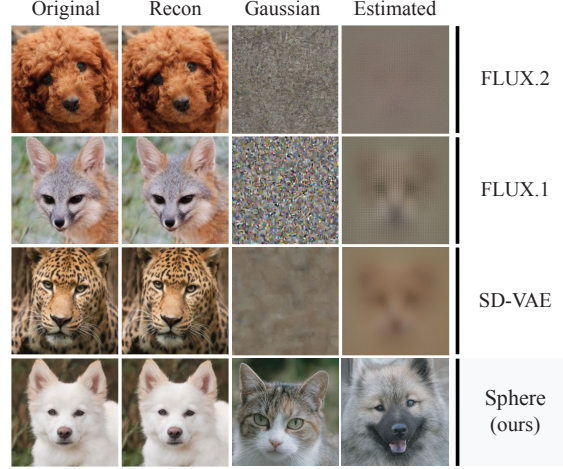


Figure 3. Posterior hole problem in VAEs. Columns: (1) Input images; (2) Autoencoder reconstructions; (3) Samples from standard Gaussian prior; and (4) Samples from estimated Gaussian posterior on Animal-Faces training set. Unlike modern FLUX.1/2 (Labs et al., 2025) and SD-VAE (Podell et al., 2024), our sphere encoder produces realistic images by decoding random points sampled from the sphere.

both uniformity and accurate reconstruction simultaneously.

Many contemporary state-of-the-art diffusion models are actually latent diffusion models (Rombach et al., 2022; Peebles & Xie, 2023; Liu et al., 2023; Ma et al., 2024; Esser et al., 2024; Podell et al., 2024; Wan et al., 2025) – hybrid models built on top of VAEs. The VAE partially Gaussianizes the image distribution, but not well enough to be sampled. A diffusion pipeline picks up the slack in the VAE, going the last mile of producing a valid latent sample for the decoder. Concurrent works have shown that more powerful representation encoders (Yu et al., 2024; Tong et al., 2026), and even spherical manifold encoders (Zheng et al., 2025), result in faster training of the diffusion layer. In our work, we show that a spherical latent space¹ can be learned so precisely that the expensive diffusion step is irrelevant.

2. Method

2.1. Spherical Latent Space

We employ an encoder E based on a Transformer (Dosovitskiy, 2020; Vaswani et al., 2017) to map an input image $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$ into a latent representation $\mathbf{z} \in \mathbb{R}^{h \times w \times d}$. The latent resolution is determined by the patch size P , such that $h = H/P$ and $w = W/P$, with d denoting the channel depth.

To construct a global spherical latent space, we define a *spherifying* function, denoted as f . This function flattens \mathbf{z} into a vector of dimension $L = h \times w \times d$ and then projects

¹In contrast to prior vMF-based approaches, we create our spherical space using simple vector RMS normalization.

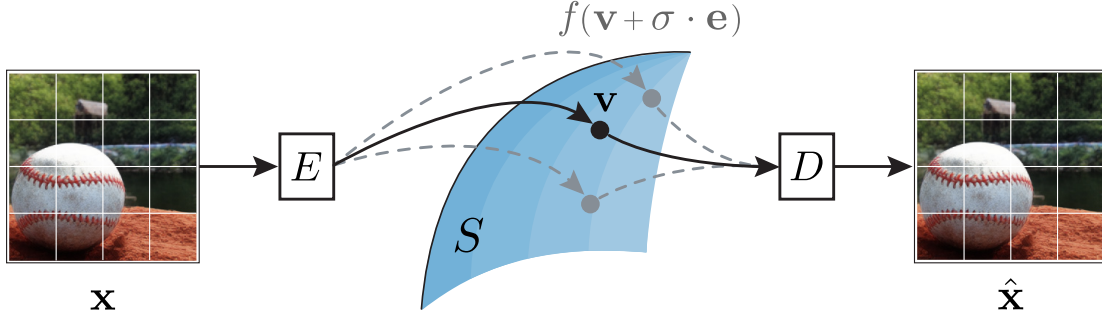


Figure 4. **Spherifying latent with noise.** Encoder E maps image \mathbf{x} to a latent, which f projects to \mathbf{v} on sphere S . During training, random Gaussian noise $\sigma \cdot \mathbf{e}$ is added to \mathbf{v} , where σ is jittered magnitude. Decoder D reconstructs the image $\hat{\mathbf{x}}$ from the re-projected noisy latent $f(\mathbf{v} + \sigma \cdot \mathbf{e})$.

it onto a sphere with radius \sqrt{L} via RMS normalization:

$$\mathbf{v} = f(\mathbf{z}) = f(E(\mathbf{x})) . \quad (1)$$

Subsequently, a decoder D reconstructs the image from \mathbf{v} :

$$\hat{\mathbf{x}} = D(\mathbf{v}) , \quad (2)$$

where $\hat{\mathbf{x}}$ denotes the reconstructed image. If the encoder maps images uniformly onto a sphere, then we can generate images by decoding random points on the sphere:

$$\hat{\mathbf{x}} = D(f(\mathbf{e})) , \quad (3)$$

where $\mathbf{e} \sim \mathcal{N}(0, \mathbf{I}) \in \mathbb{R}^L$ is random anisotropic Gaussian and $f(\mathbf{e})$ is uniformly distributed on the sphere. For simplicity, we use $\hat{\mathbf{x}}$ to denote the decoder output in both reconstruction and generation scenarios.

2.2. Spherifying with Noise

Our training process uses embedding vectors of natural images, and also noisy versions of those embedding vectors. The purpose of training with noisy vectors is two-fold. First, noisy clouds of vectors densely cover the latent space, enabling us to train the decoder on the continuous global latent sphere, rather than only on the finite set of embedding vectors. Second, by using a loss that promotes accurate decoding of noisy latent vectors, we force the noisy clouds produced by each training image to spread apart and cover the entire latent sphere.

From a Normal distribution, we randomly sample a noise vector $\mathbf{e} \sim \mathcal{N}(0, \mathbf{I}) \in \mathbb{R}^L$ to perturb the direction of \mathbf{v} :

$$\mathbf{v}_{\text{NOISY}} = f(\mathbf{v} + \sigma \cdot \mathbf{e}) , \quad (4)$$

where the scalar σ controls the noise magnitude. Note that f is applied again here to project the perturbed vector back onto the spherical surface.

Jittering Sigma. To cover diverse directions on the sphere, we jitter σ during the training. By sampling a scalar r

uniformly from $[0, 1]$, we compute σ as:

$$\sigma = r \cdot \sigma_{\max} , \quad (5)$$

where the σ_{\max} is the maximum noise limit. The case of $r = 0$ reduces to the naive spherifying in Equation (1). Later we determine the optimal value for σ_{\max} with experiments. This core design is illustrated in Figure 4.

2.3. Training Objective

Consider two perturbed latent vectors, $\mathbf{v}_{\text{NOISY}}$ and $\mathbf{v}_{\text{noisy}}$, with large and small noise. $\mathbf{v}_{\text{NOISY}}$ is defined as in Equation (4) with $\sigma \in [0, \sigma_{\max}]$. The other perturbed $\mathbf{v}_{\text{noisy}}$ has less jitter:

$$\mathbf{v}_{\text{noisy}} = f(\mathbf{v} + \sigma_{\text{sub}} \cdot \mathbf{e}) , \quad (6)$$

where $\sigma_{\text{sub}} = s \cdot \sigma$, and s is uniformly sampled from $[0, 0.5]$. Note that $\mathbf{v}_{\text{noisy}}$ shares the same noise direction \mathbf{e} as $\mathbf{v}_{\text{NOISY}}$.

Pixel Reconstruction Loss. This loss ensures that the decoder is an approximate inverse of the encoder, and that the decoder creates valid images. We have the standard pixel-level reconstruction loss, which combines of smoothed L1 loss (Girshick, 2015) and perceptual loss (Johnson et al., 2016). This loss encourages the decoder to reconstruct the input image \mathbf{x} from its noisy latent representation $\mathbf{v}_{\text{noisy}}$:

$$\mathcal{L}_{\text{pix-recon}} = \mathcal{L}_{\text{L1} + \text{perceptual}}(D(\mathbf{v}_{\text{noisy}}), \mathbf{x}) . \quad (7)$$

Pixel Consistency Loss. This consistency loss ensures that the latent space is smooth and well structured by promoting that nearby latent vectors produce similar images:

$$\mathcal{L}_{\text{pix-con}} = \mathcal{L}_{\text{L1} + \text{perceptual}}(D(\mathbf{v}_{\text{NOISY}}), \text{sg}(D(\mathbf{v}_{\text{noisy}}))) , \quad (8)$$

which also uses the combination of smooth L1 loss and perceptual loss, and $\text{sg}(\cdot)$ denotes stop-gradient operation.

Latent Consistency Loss. It is well known that image similarity is better measured in latent space than in pixel space (Zhang et al., 2018; Radford et al., 2021). This is

the reason why our pixel similarities use a perceptual loss, which relies on features produced by a static VGG model. To achieve a stronger consistency loss, we also measure image similarity using the latent space of our own encoder.

We want a natural image \mathbf{x} and its noisy decoded representation $D(\mathbf{v}_{\text{NOISY}})$ to be semantically similar. The semantic similarity is measured by applying the encoder to both, and computing the cosine similarity between their latent representations. This yields the following loss:

$$\mathcal{L}_{\text{lat-con}} = \mathcal{L}_{\text{cosine similarity}}(\mathbf{v}, E(D(\mathbf{v}_{\text{NOISY}}))) \quad (9)$$

This loss serves an additional important purpose: It improves the iterative generation process we discuss later by encouraging the encoder to map distorted images, $D(\mathbf{v}_{\text{NOISY}})$, that may be off the image manifold to “cleaned up” latent vectors that reflect on-manifold images.

Overall Loss. The overall training loss is a weighted sum of the three components:

$$\mathcal{L} = \mathcal{L}_{\text{pix-recon}} + \mathcal{L}_{\text{pix-con}} + \mathcal{L}_{\text{lat-con}} \quad (10)$$

More details about loss weights and training hyperparameters are provided in Appendix D.

2.4. Model Architecture

Our architecture employs the standard ViT (Dosovitskiy, 2020) for both encoder and decoder. We insert 4-layer MLP-Mixers (Tolstikhin et al., 2021) in the end of the encoder (before spherification) and the beginning of the decoder. This aims to improve cross-token mixing and globalization of features without the expense of linear layers on the full flattened vector. A final RMSNorm layer (Zhang & Sennrich, 2019) with learned affine parameters is added to each MLP-Mixer to bound the latent magnitude ($\leq \sqrt{L}$). This regularization proves critical for stabilizing training, especially when there is a dramatic divergence between the decoder outputs of $\mathbf{v}_{\text{noisy}}$ and $\mathbf{v}_{\text{NOISY}}$. We use both RoPE (Su et al., 2024) positional embedding and sinusoidal absolute positional encoding. We found that removing the sinusoidal positional embedding hurts generation quality.

For class-conditional generation, we implement AdaLN-Zero (Perez et al., 2018; Peebles & Xie, 2023) in both the encoder and decoder using separate class embeddings. A learned `null` embedding is trained for classifier-free guidance (CFG) (Ho & Salimans, 2022) with a probability of 0.1. For image reconstruction tasks, we found using either random class embeddings or the `null` embedding is effective in the conditional setting. We default to the `null` embedding for simplicity. In addition, CFG can be applied in either the latent space (after the encoder), or the pixel space (after the decoder), or both.

Algorithm 1 summarizes the forward pass for one-step or few-step generation at inference time. We count $D(f(\mathbf{e}))$ as

Algorithm 1 pseudocode of generation forward pass.

```
e = Normal(0, 1).sample([L]) # random noise vector

# one-step generation
v = spherify(e, sampling=False) # Eq.(1) w/o noise
x = D(v, y) # y is class embedding

if do_dec_cfg: # in pixel space
    x_uncond = D(v, y_null) # null class embedding
    x = x_uncond + cfg * (x - x_uncond)

# few-step refinement
for _ in range(T - 1): # T steps in total
    z = E(x, y)

    if do_enc_cfg: # in latent space
        z_uncond = E(x, y_null)
        z = z_uncond + cfg * (z - z_uncond)

    v = spherify(z, sampling=True) # Eq.(4) w/ noise
    x = D(v, y)

    if do_dec_cfg:
        x_uncond = D(v, y_null)
        x = x_uncond + cfg * (x - x_uncond)

return x
```

one-step generation, and few-step generation as iteratively encoding and decoding $T - 1$ times². We fix $r = 1.0$ in Equation (5) to use the maximum noise magnitude across all steps. For reconstruction task, no noise is added ($r = 0.0$).

3. Quantitative Experiments

We adopt generation FID (gFID) (Heusel et al., 2017) and Inception Score (IS) (Salimans et al., 2016) to assess generation quality, while reconstruction FID (rFID) measures reconstruction quality. All metrics are computed using 50K randomly sampled training images. For class-conditional generation, we have a balanced distribution with an equal number of random samples per class.

We perform experiments on CIFAR-10 (Krizhevsky et al., 2009) with small image size 32×32 , as well as ImageNet (Russakovsky et al., 2015), Animal-Faces (Choi et al., 2020), and Oxford-Flowers (Nilsback & Zisserman, 2008) with large image size 256×256 . Center cropping and horizontal flipping with 0.5 probability are the only data augmentation.

3.1. Small Image Size

We first comprehensively test our method with small image size (32×32) on CIFAR-10 (Krizhevsky et al., 2009). We build encoder and decoder with ViT Large, which consists of 24 layers and yields a latent dimension $L = 16 \times 16 \times 8$. The model, indicated as Sphere-L, is trained for 5000 epochs for conditional generation, 10000 epochs for unconditional generation, following the other setup in Appendix D.

²While a ‘step’ represents a single model iteration regardless of CFG, NFE (number of function evaluation) counts the dual forward passes required by CFG.

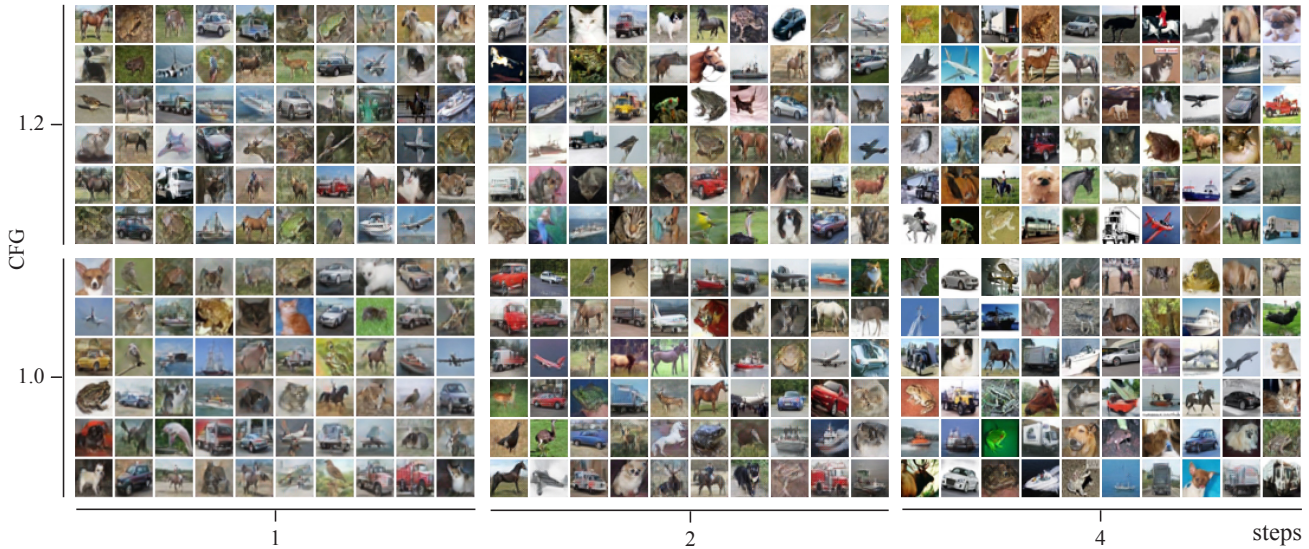


Figure 5. **Uncurated CIFAR-10 conditional generation** with different sampling steps and with/without CFG. Convincing images can be formed with a single forward pass, with reliability and gFID improving with up to 4 steps.

Table 1. **Few-step generation results** on CIFAR-10.

METHOD	STEPS	RFID ↓	gFID ↓	IS ↑
CONDITIONAL GENERATION W/O CFG				
SPHERE-L	1	0.59	18.68	9.1
	2	-	8.28	9.9
	4	-	2.72	10.5
	6	-	1.65	10.7
STYLEGAN2 (KARRAS ET AL., 2020)	1	-	6.96	9.53
STYLEGAN2 + ADA (KARRAS ET AL., 2020)	1	-	3.49	10.2
UNCONDITIONAL GENERATION				
SPHERE-L	1	0.53	35.67	6.7
	2	-	14.13	8.4
	4	-	4.31	9.8
	6	-	2.34	10.2
DDIM (SONG ET AL., 2021)	1K	-	3.17	-
DDPM (HO ET AL., 2020)	1K	-	3.17	9.4
IMPROVED-DDPM (NICHOL ET AL., 2021)	4K	-	2.90	-

Table 1 presents both conditional and unconditional generation results. For conditional generation, our method achieves strong results in both one-step and few-step generation, even without CFG. For example, with less than 10 sampling steps, sphere encoder yields gFID way below 2.0 and IS above 10. For unconditional generation task, our method achieves better gFID and IS with less than 10 sampling steps, a $100\times$ reduction in sampling steps, comparing to diffusion-based methods (Ho et al., 2020). Figure 5 depicts qualitative results with different steps and CFG. Visually, our generation results without CFG look the same as those using CFG. Appendix A provides quantitative results on CIFAR-10 with CFG. Additionally, we discuss the memorization risk when training on small datasets like CIFAR-10 with extensive epochs in Appendix B.

Table 2. **Few-step generation results** (gFID ↓) on Animal-Faces (AF) and Oxford-Flowers (OF).

SPHERE-L MODEL \ STEPS	1	2	4	6
OF W/ CFG=1.6	25.12	14.08	11.25	10.63
OF	27.12	16.60	12.98	12.26
AF	21.70	19.29	18.23	17.97
STYLEGAN2 (KARRAS ET AL., 2020)				
AF - CAT	5.13	-	-	-
AF - DOG	19.37	-	-	-
AF - WILD	3.48	-	-	-

3.2. Large Image Size

We then evaluate our method with large image size (256×256) on Oxford-Flowers (Nilsback & Zisserman, 2008) (8K images), Animal-Faces (Choi et al., 2020) (16K images), and ImageNet (Russakovsky et al., 2015) (1.2M images).

Animal-Faces and Oxford-Flowers. We employ a ViT Large for encoder and decoder, *i.e.*, Sphere-L, with a latent dimension of $L = 32 \times 32 \times 128$, training for 1000 epochs. Due to the low diversity of Animal-Faces (3 classes), we train unconditional models, while for Oxford-Flowers, we utilize conditional generation for the 102 classes. For evaluation, we adhere to the standard protocol of randomly sampling 50K images, even though the training sets are relatively small. Table 2 reports quantitative results on both datasets. We report metrics only up to 6 sampling steps, as performance saturates beyond this point. Uncurated qualitative results are provided in Figures 17 and 18.

ImageNet. We train class-conditional ViT-Large/-XLarge models on ImageNet for 800 epochs, utilizing a latent dimension of $L = 32 \times 32 \times 64$. Table 3 evaluates our approach alongside GANs, diffusion models, and other direct pixel-space generation frameworks. At comparable

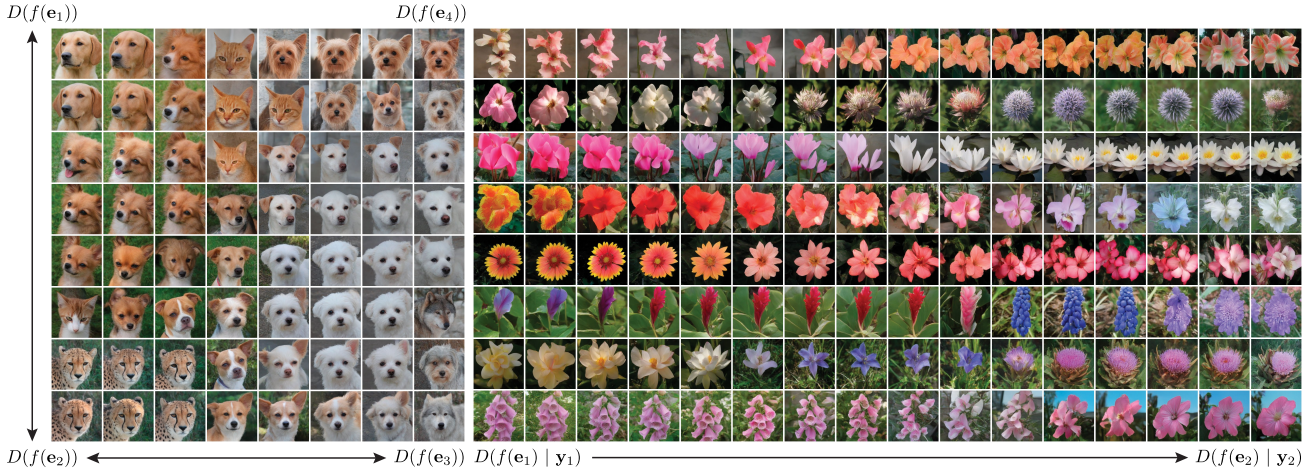


Figure 6. **Latent interpolation** on Animal-Faces and Oxford-Flowers. Images are generated in 4 steps without CFG. (left) Interpolation in a 2D space that spans 4 synthetic images. (right) Each row interpolates between a random vector e on the sphere and a class conditional vector y . Note that our model exhibits fast/sudden transitions between image classes rather than producing “hybrid” images that unrealistically merge properties of different object types. This property is necessary for a model to reliably convert random samples from the sphere into realistic images, as it makes the probability of observing a hybrid image small.

parameter counts, our sphere encoder achieves competitive FID scores while requiring fewer sampling steps; its performance falls within the range of recent high-performing models, outperforming several established baselines.

3.3. Lower FID scores?

Because low FID scores do not always align with perceptual realism (Stein et al., 2023), Table 3 reports FIDs for our qualitatively strongest models, prioritizing visual quality over the optimization of a single numerical metric. Lower FID scores are possible with some tradeoffs. For example, while training on CIFAR-10 for 10K epochs can reduce FID to 0.94, it does not yield a proportional gain in visual clarity. On ImageNet, increasing sampling steps beyond 4 improves the FID to 3.9 and sharpens local edges, yet this numerical gain can introduce more abstract object structures. This phenomenon – where FID rewards local texture refinement even at the cost of global semantic coherence (Jayasumana et al., 2024) — highlights a nuanced trade-off in generative modeling that warrants further investigation.

Still, our reported FID scores trail some recent high-performance generative models. We suspect this may be due to our use of pixel-space similarity losses, which likely contributes to the subtle edge blurring observed in our uncured results (Figures 15 and 16). In contrast, other recent works achieve high sharpness and low-FID through similarity metrics based purely on latent-space representations or multi-stage GAN losses – a direction that should be evaluated in future work.

Finally, note that our training and conditioning methods do not rely on the discreteness of the ImageNet ontology. This generality opens the door for our methods to be transferred to the text-to-image setting.

Table 3. **Few-step generation results** on ImageNet 256×256 .

	METHOD	PARAMS	STEPS	GFID ↓	IS ↑
GANS	BIGGAN-DEEP (b21)	56M	1	6.90	171.4
	STYLEGAN-XL (s22)	166M	1	2.30	265.1
	GIGAGAN (k23)	569M	1	3.45	225.5
DIFFUSIONS	ADM-G (d21)	554M	250	4.59	186.7
	ADM-G	554M	25	5.44	-
	SID (h23)	2B	1000	2.44	256.3
	SID2 (h24)	-	512	1.38	-
	JIT-H (l25)	953M	100	1.86	303.4
	JIT-G	2B	100	1.82	292.6
OURS	JETFORMER (t25)	2.8B	-	6.64	-
	FRACTALMAR-H (l25)	848M	96	6.15	348.9
	SPHERE-L	950M	4	4.76	301.8
	SPHERE-XL	1.3B	4	4.02	265.9

4. Qualitative Experiments

Latent Interpolation. In Figure 6, we interpolate between latent vectors on the Animal-Faces and Oxford-Flowers models to investigate the learned latent manifold. On Animal-Faces, we randomly sample four noise vectors e in Equation (3) and visualize their corresponding images at the corners of the figure. We interpolate the latent space bilinearly to fill in the other images. On Oxford-Flowers, we randomly sample a noise vector e and a class for each side of each row. Since the model is class-conditional, we interpolate both input noise and class embeddings linearly as we move horizontally across each row.

We see that the model exhibits *fast transitions* between object types as we move through latent space. For example, starting with the bottom-left image of a cheetah, we observe a sudden transition from cheetah to cat as we move vertically, and from cheetah to dog as we move horizontally. The model does not linger in a half-cheetah / half-dog state that is absent in the training data. These fast-transitions are

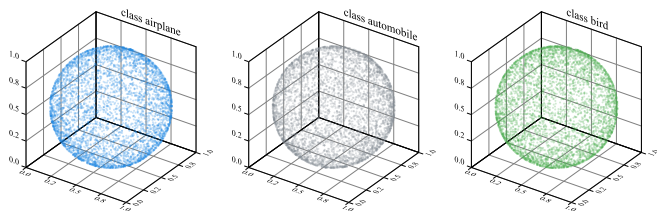


Figure 7. **Latent space visualization** using random projection on CIFAR-10 training set. Each sphere shows the latent vectors of a different class. The conditional latent distributions appear approximately uniform.

required of a reliable generative model, as the probability of sampling an impossible/hybrid image should be minimized. This important property of the sphere encoder differentiates it from other latent models. GANs, for example, tend to exhibit slow transitions, resulting in frequent production of distorted objects, *e.g.*, Figure 8 and 9 in (Brock et al., 2019).

Conditional Uniformity. The latent distribution of our model should be uniform on the sphere. For conditional models, the latent distribution must be *conditionally* uniform. To understand why, consider a conditional model with two classes, cat and dog. Suppose we have an unconditional encoder and a conditional decoder. It is likely that an unconditional encoder would structure latent space with cats in one region and dogs in another. Even if the union of the two classes achieves uniform coverage, conditional generation may fail; we cannot reliably decode a dog from a random vector, as half the time this vector will represent a cat.

We avoid this pitfall by using a conditional encoder. In this case, our training objective naturally creates a latent distribution that is conditionally uniform, *i.e.*, dogs alone uniformly

cover the sphere, as do cats. As a result, any uniformly sampled vector can be decoded to create the desired object.

To better understand the latent space learned by our method, we visualize it in Figure 7. We extract \mathbf{v} in Equation (1) for all CIFAR-10 50K training samples. We project latents to 3D space using a random Gaussian matrix, and then normalize each projected vector to have unit length.

We visualize the results separately for three random classes (other classes look similar). We see that the latent space achieves conditional uniformity – we get even coverage of the sphere for each class in isolation.

5. Image Editing

This section presents two *training-free* editing applications that leverage the expressivity and robustness of our latent space: simple semantic manipulation and image crossover.

Conditional Manipulation. Given an out-of-domain image, *e.g.*, a “woolly panda” in Figure 8, we can manipulate it by conditioning on different ImageNet classes. We simply encode and decode the image with multiple steps, using the Sphere-L model (in Table 3) trained on ImageNet. We set a fixed noise strength $r = 1.0$ in Equation (5) and $\gamma = 0$ in Equation (13) across all steps, and do not apply CFG.

Figure 8 demonstrates that a single step effectively captures the primary structure of the input while adapting its texture to align with the target class. Subsequent iterations (*e.g.*, 4-step generation) further refine these class-specific characteristics and textures, achieving semantic alignment while preserving the original image’s structural integrity.

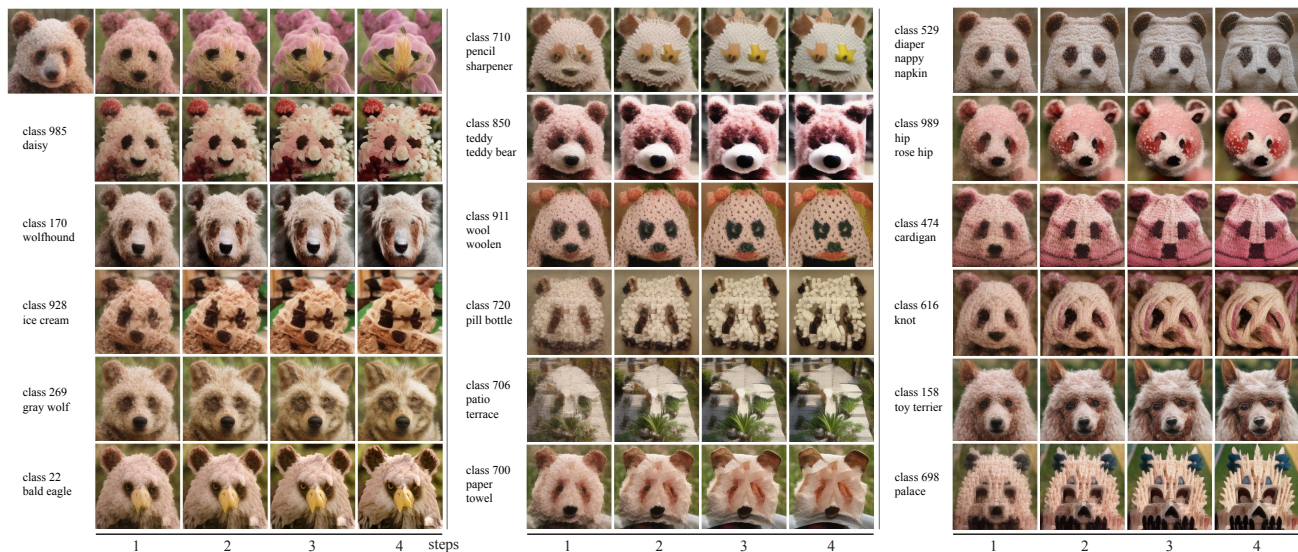


Figure 8. **Conditional manipulation** via iterative encoding and decoding on ImageNet model. We demonstrate the model’s expressivity using an out-of-domain input (a “woolly panda” top-left). Each row shows the result of conditioning the iterative process on different ImageNet classes without CFG. For example, the first row is conditioned on class 580 (greenhouse, nursery, glasshouse).

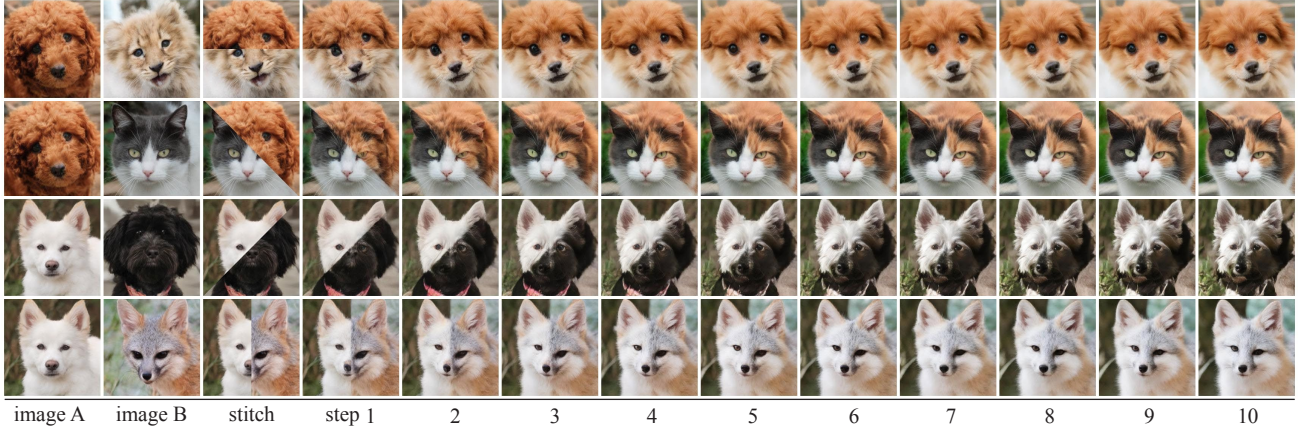


Figure 9. **Image crossover** using the sphere encoder trained on Animal-Faces. A composite of two images (A and B) is iteratively processed through the encoder-decoder pipeline until it converges to a coherent sample on the learned image manifold.

Image Crossover. We further demonstrate the model’s capability for “image crossover” by processing manually stitched composites of distinct source images (Figure 9). This process similarly operates without CFG. By repeatedly encoding and decoding the stitched composite (up to 10 steps), the model naturally harmonizes the content and smooths boundary discontinuities. For these experiments, we set the noise magnitude $r = 0.25$ (Equation (5)) and apply a decay schedule (Equation (13)) with $\gamma = 1$, which we found yields the most seamless blending.

This iterative refinement forces the manipulated image to converge toward a valid point on the learned spherical manifold. Notably, unlike diffusion models, *e.g.*, Figure 12 in (Liu et al., 2023), that require noise injection before projecting onto the image manifold, our encoder directly projects the stitched image into the latent space without adding noise (through the encoder). This deterministic path preserves the semantic integrity of the original sources while ensuring a fluid, natural transition between them.

6. Main Ablations

This section presents ablation studies on key design choices of our method. Additional ablations regarding CFG strategies, noise distribution, uniform regularization, volume compression ratio, and model size are provided in Appendix C.

Determining Noise Magnitude. In this section, we analyze the maximum noise magnitude σ_{\max} in Equation (4) from a geometric perspective to understand its impact, and empirically determine its optimal value. We begin with the noise-to-signal ratio (NSR) η as the ratio of the expected noise magnitude to the signal magnitude in high-dimensional space:

$$\eta = \mathbb{E} \frac{\|\mathbf{e}\|}{\|\mathbf{v}\|} = \frac{\sigma_{\max} \sqrt{L}}{\sqrt{L}} = \sigma_{\max}. \quad (11)$$

Because of the concentration of measure phenomenon, \mathbf{e} is nearly orthogonal to \mathbf{v} , *i.e.*, $\mathbf{v}^\top \mathbf{e} = 0$ when $L \rightarrow \infty$. Thus, the angle α formed between $\mathbf{v} + \sigma \cdot \mathbf{e}$ and \mathbf{v} satisfies:

$$\tan(\alpha) \approx \frac{\|\mathbf{e}\|}{\|\mathbf{v}\|} \approx \sigma_{\max} = \eta, \quad (12)$$

which gives σ_{\max} an interpretable geometric meaning. The noise magnitude σ_{\max} can be equivalently expressed by either the angle α or the NSR η .

We build a conditional encoder and decoder using classic ViT with Base size (12 layers), and train them on ImageNet for 200 epochs. The latent dimension is $L = 16 \times 16 \times 256$. We vary α from 45° to 88° , corresponding to $\sigma_{\max} = \tan(\alpha)$ from 1 to 28.6.

Our first ablations are done to select the noise level. The NSR η guides the difficulty of the decoder’s task. The decoder aims to generate the same image from $\mathbf{v}_{\text{NOISY}}$ as from the clean \mathbf{v} . When $\alpha \leq 45^\circ$ (equivalently $\eta \leq 1$), the noise does not overwhelm \mathbf{v} and the latent clouds generated by each training point are well separated. The decoder reconstructs images well, but the noisy latents fail to cover the sphere. In this case, generation from random sampling using Equation (3) fails, with a high gFID in Figure 10.

As $\alpha \rightarrow 90^\circ$, the latent representations of images are forced apart and the decoder starts to generate images from random latents. The gFID drops significantly in Figure 10. Figure 11 shows some sampled images with different α , demonstrating that a lower α leads to abstract and blurry images, while a higher α produces more realistic details.

Once we dial in α , we find that we can fix it for various latent dimensions L because of the dimension-invariant property of angle α in Equation (11). However, we found that the optimal α varies slightly with image size. For small image size, *e.g.*, 32×32 on CIFAR-10, we found $\alpha = 80^\circ$ works best. For large image size, *e.g.*, 256×256 on ImageNet, we found $\alpha = 85^\circ$ works best.

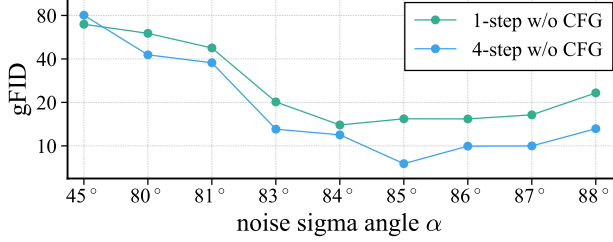
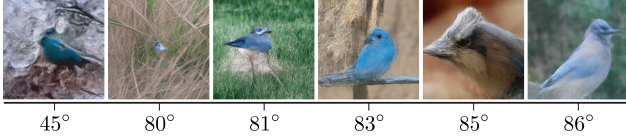

 Figure 10. Quantitative impact of the angle α on ImageNet.

 Figure 11. Qualitative impact of the angle α on ImageNet with 4-step generation.

Table 4. Ablation of loss terms on ImageNet. Loss terms are added incrementally from top to bottom.

LOSS \mathcal{L}	STEPS	RFID ↓	W/O CFG		W/ CFG = 1.6	
			GFIG ↓	IS ↑	GFIG ↓	IS ↑
$\mathcal{L}_{\text{PIX-RECON}}$	1	1.70	25.35	168.8	20.68	204.4
	4	-	13.58	162.9	11.76	231.7
+ $\mathcal{L}_{\text{PIX-CON}}$	1	0.89	17.19	183.8	14.74	218.3
	4	-	10.12	189.9	10.47	247.1
+ $\mathcal{L}_{\text{LAT-CON}}$	1	1.32	15.40	188.3	13.14	225.1
	4	-	7.53	176.4	7.28	243.1

Training Loss. We evaluate the individual contribution of each proposed loss term in Section 2.3 to generation quality. Adopting the same ImageNet setup as the previous ablation, Table 4 reports the results of adding each term incrementally.

Starting with only the pixel reconstruction loss $\mathcal{L}_{\text{pix-recon}}$, we observe the model can generate images, but the quality is suboptimal with a serious “waffle” artifact. Adding the pixel consistency loss significantly improves generation quality by removing the artifact, as it encourages the decoder to produce consistent images from perturbed latents.

Including the latent consistency loss yields the best performance, demonstrating its effectiveness in guiding the encoder-decoder pair to learn a coherent latent manifold. Figure 12 visualizes the optimization path from noisy latent to clean latent on the sphere for those two consistency losses during training. Overall, each loss contributes positively to the model’s ability to generate high-quality images.

Latent Spatial Resolution. In these ablations, we keep the latent dimension constant and vary the latent spatial resolution by adjusting channel depth d of latent dim L . Table 5 presents the results on ImageNet, suggesting that a higher latent spatial resolution with volume compression ratio 3.0 works best on ImageNet, *i.e.*, $L = 32^2 \times 64$. On CIFAR-10, Animal-Faces, and Oxford-Flowers, we also observed that a higher latent spatial resolution yields better

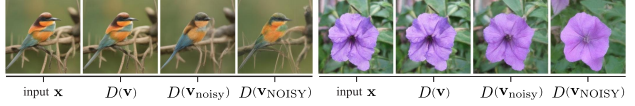


Figure 12. Consistency optimization path from noisy latent to clean latent on the sphere, pushing the decoder to generate consistent and diverse images from right to left.

Table 5. Ablation on latent spatial resolution with two optimal compression ratios on ImageNet.

LATENT L	STEPS	RFID \downarrow	W/O CFG		W/ CFG = 1.2	
			GFIG \downarrow	IS \uparrow	GFIG \downarrow	IS \uparrow
VOLUME COMPRESSION RATIO = 1.5						
$16^2 \times 512$	1	0.64	16.06	169.8	14.25	187.6
	4	-	10.03	153.0	8.71	180.8
$32^2 \times 128$	1	0.60	20.63	145.1	17.69	164.9
	4	-	15.45	127.0	12.75	154.6
VOLUME COMPRESSION RATIO = 3.0						
$16^2 \times 256$	1	0.74	16.93	206.3	15.39	227.7
	4	-	8.30	202.9	7.85	232.3
$32^2 \times 64$	1	0.61	16.80	227.8	15.33	267.0
	4	-	7.65	219.4	7.48	252.4

generation quality, but with a lower compression ratio 1.5. The finalized latent dimensions are detailed in Appendix D.

Sampling Schemes. In Algorithm 1, the few-step generation involves adding noise in spherifying at each step following Equation (4). We investigate two key aspects of this noise injection mechanism: strength schedule and sampling schedule. First, regarding the noise strength controlled by r in Equation (5), we compare two cases: (a) a fixed $r = 1.0$ for all steps, and (b) a decaying schedule, where r decreases from 1.0 to a minimum value. Second, we also have two options for sampling the noise vector e in Equation (4): (a) sampling independent noise e for each step, versus (b) sharing the same noise e across all steps. We qualitatively and quantitatively evaluate both aspects without CFG.

We simply decay the noise strength r with a linear schedule:

$$r = \left(1 - \frac{t-1}{T-1}\right)^\gamma, \quad (13)$$

where t is the current step from 2 to T in the loop, and decay rate $\gamma = 1$. The $t = 1$ is the first step of decoder forward pass, which is not involved with spherifying process. When $\gamma = 0$, it corresponds to the fixed schedule, *i.e.*, $r = 1.0$.

Table 6 presents the results on ImageNet with those four sampling schemes. The fixed schedule ($\gamma = 0$) outperforms the decaying schedule ($\gamma = 1$) across all sampling steps. In addition, sharing the same noise e across all steps consistently yields better results than using independent noise for each step. We hypothesize that sharing the same noise helps maintain a consistent direction during the optimization path on the sphere, leading to more stable and effective generation. Overall, the best sampling scheme is using a fixed

Table 6. **Ablation on sampling schemes** with few-step generation and no CFG. Results are reported with gFID \downarrow on ImageNet.

γ	SHARE \mathbf{e}	1	2	4	6	8	10
0	-	16.57	7.99	7.68	7.49	7.61	7.71
	✓	16.66	7.97	5.99	7.78	10.11	12.25
1	-	16.73	8.12	9.02	12.17	15.08	18.09
	✓	16.81	8.02	8.53	12.32	17.50	22.78

noise strength with shared noise across all steps.

Figure 13 shows generated images on both CIFAR-10 and ImageNet under different sampling schemes. Visual inspection confirms that sharing the same noise across steps yields significantly more coherent and higher-quality results than using independent noise. Notably, on ImageNet, we observe that sharing noise with decay schedule $\gamma = 1$ produces exceptionally sharp images as the number of sampling steps increases. For example, with 10 steps, the images exhibit a distinct, hyper-sharp aesthetic reminiscent of paper art.

Achieving a Uniform Distribution. Our method does not employ explicit regularization to achieve a uniform latent distribution on the sphere. Comprehensive ablation studies in Appendix C.3 demonstrate that such regularization is unnecessary, as our method naturally achieves the desired latent properties.

7. Related Work

Spherical latent space has been explored primarily through variational inference, using priors such as the von Mises-Fisher distribution (Xu & Durrett, 2018; Davidson et al., 2018; De Cao & Aziz, 2020; Ke & Xue, 2025). However, these approaches inherit limitations from VAEs, including significant posterior-prior mismatch. These issues are compounded in high-dimensional latent spaces, where sampling relies on intricate reparameterization or rejection sampling techniques. Although (Zhao et al., 2019) drew inspiration from StyleGAN (Karras et al., 2019; 2020; 2018) to enable direct sampling in a high-dimensional unit spherical space via simple normalization, their experiments were limited to toy datasets like MNIST. Since VAEs can already perform direct sampling on MNIST (Dai & Wipf, 2019; Rezende & Viola, 2018), the advantages of spherical latent spaces in this context remain unclear. Crucially, the potential of high-dimensional spherical latent spaces (Kumar & Patel, 2026) for image generation remains significantly underexplored.

Few-step generation has been extensively studied in both GANs and diffusion models. GANs (Goodfellow et al., 2014) are inherently created for one-step generation. While approaches like ProgressiveGAN (Karras et al., 2018) introduce multi-stage refinement, these primarily serve as training strategies (Denton et al., 2015; Zhang et al., 2017; Brock et al., 2019) rather than altering the fundamental single-pass

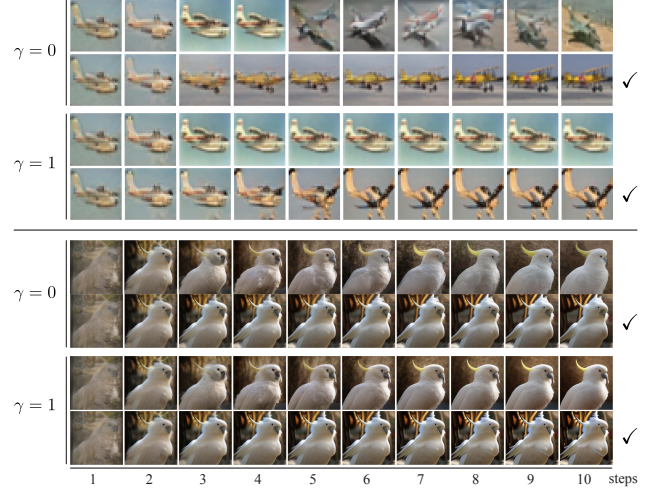


Figure 13. **Qualitative impact of sampling schemes** on generation without CFG for CIFAR-10 (top) and ImageNet (bottom). The ✓ indicates sharing the same noise \mathbf{e} across steps. Shared noise with $\gamma = 1$ yields superior coherence and a sharp “paper art” aesthetic on ImageNet as sampling steps increase.

inference process. Conversely, diffusion models rely on an iterative generation process, typically requiring hundreds to thousands of steps. Although distillation (Frans et al., 2024; Yin et al., 2024; Xie et al., 2024; Salimans & Ho, 2022; Geng et al., 2025; Sauer et al., 2024) and consistency techniques (Song et al., 2023; Geng et al., 2024; Yang et al., 2024; Lu & Song, 2025) have been proposed to accelerate this to a few steps, their core insight aims to approximate the original diffusion trajectory.

Pixel-space generation is common for GANs (Karras et al., 2020; Kang et al., 2023; Sauer et al., 2022; Brock et al., 2019; 2018) but challenging for diffusion models due to the high dimensionality of pixel space (Child, 2019; Van den Oord et al., 2016; Chen et al., 2020; Hawthorne et al., 2022; Yu et al., 2023). Recent advances based on diffusion mechanisms have made strides in pixel-space generation (Li et al., 2025; Li & He, 2025; Yu et al., 2025). Our sphere encoder diverges fundamentally from both paradigms. Its few-step mechanism iteratively traverses between latent and pixel spaces, grounded in a spherical latent space.

Signal processing. Concepts from this work take inspiration from sphere encoders/decoders in wireless networking that distribute codewords uniformly across a sphere (Studer et al., 2008; Studer & Bölcskei, 2010).

8. Conclusion

The Sphere encoder is a novel generative framework that enables few-step, high-fidelity image synthesis. Our key observation is that the distribution of latents can be tightly controlled on a uniform sphere, enabling the training of an autoencoder that can be directly sampled.

This work is intended to be a proof-of-concept for direct conditional and unconditional generation from an autoencoder, and we suspect this first implementation is far from optimal. To illustrate the average generation quality, Figure 14 shows randomly selected ImageNet results, with comprehensive uncured examples provided in Figures 15 and 16.

One drawback of our approach is that parameters must be allocated for both the encoder and decoder, and two passes are needed through the encoder during training, *i.e.*, one for latent encoding and another for consistency loss. Interestingly, improvements that enable single-pass generation for more complex distributions would eliminate the need for the encoder at inference time, and possibly at training time as well. We think there are many avenues for future research that could unlock this and other capabilities, *e.g.*, text-to-image generation.

Acknowledgements

We would like to thank Tan Wang, Chenyang Zhang, Tian Xie, Wei Liu, Felix Juefei Xu, and Andrej Risteski for their valuable discussions and feedback.

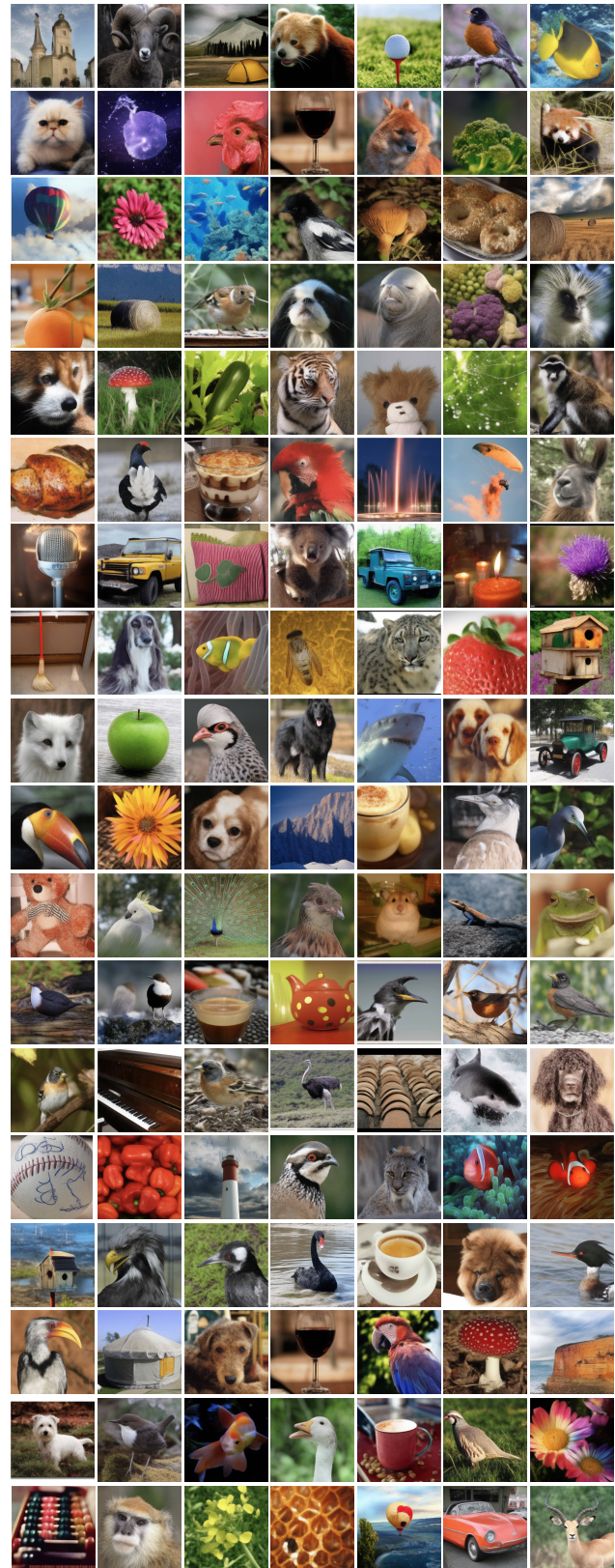


Figure 14. Randomly selected images generated by the Sphere encoder for ImageNet (256×256). Results are generated using Sphere-XL with 4-step sampling and CFG = 1.4.



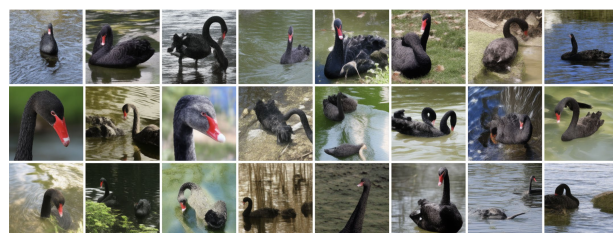
class 9: ostrich, struthio camelus



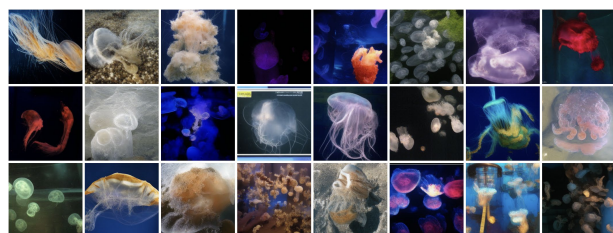
class 17: jay



class 42: agama



class 100: black swan, cygnus atratus



class 107: jellyfish



class 130: flamingo



class 143: oystercatcher, oyster catcher



class 145: king penguin, aptenodytes patagonica



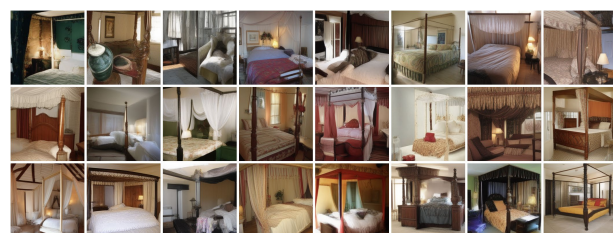
class 279: arctic fox, white fox, alopex lagopus



class 387: lesser panda, red panda



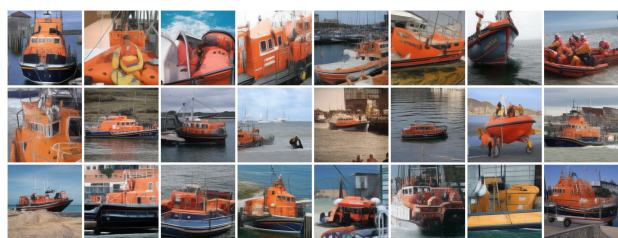
class 425: barn



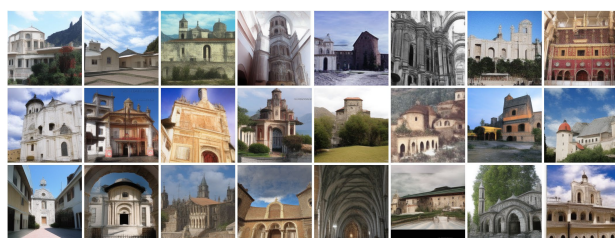
class 564: four-poster

Figure 15. **Uncurated results** on ImageNet (256×256). Results are generated using Sphere-XL with 4-step sampling and $\text{CFG} = 1.4$.

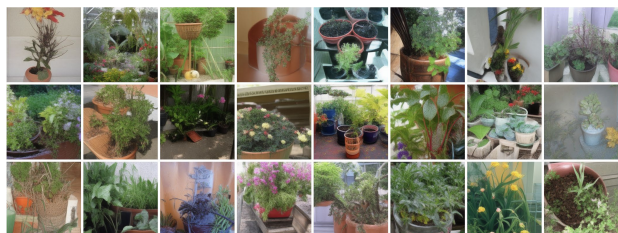
Image Generation with a Sphere Encoder



class 625: lifeboat



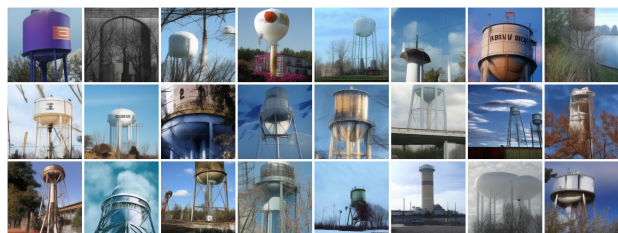
class 663: monastery



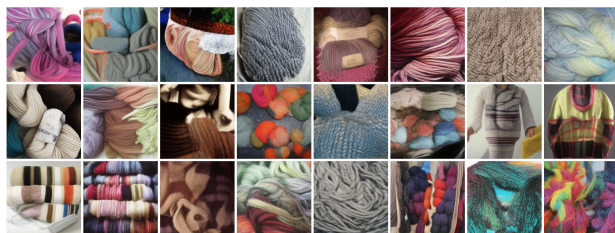
class 738: pot, flowerpot



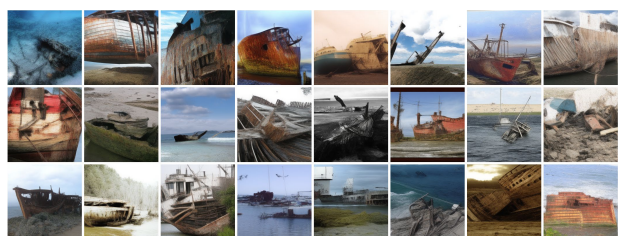
class 850: teddy, teddy bear



class 900: water tower



class 911: wool, woolen, woollen



class 913: wreck



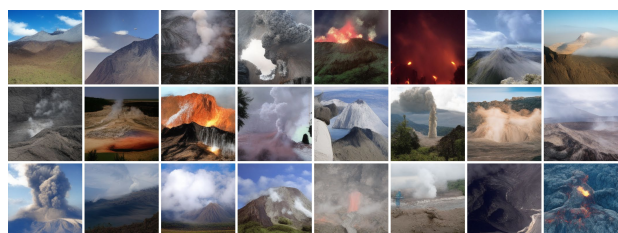
class 927: trifle



class 946: cardoon



class 957: pomegranate



class 980: volcano



class 985: daisy

Figure 16. **Uncurated results** on ImageNet (256×256). Results are generated using Sphere-XL with 4-step sampling and $\text{CFG} = 1.4$.

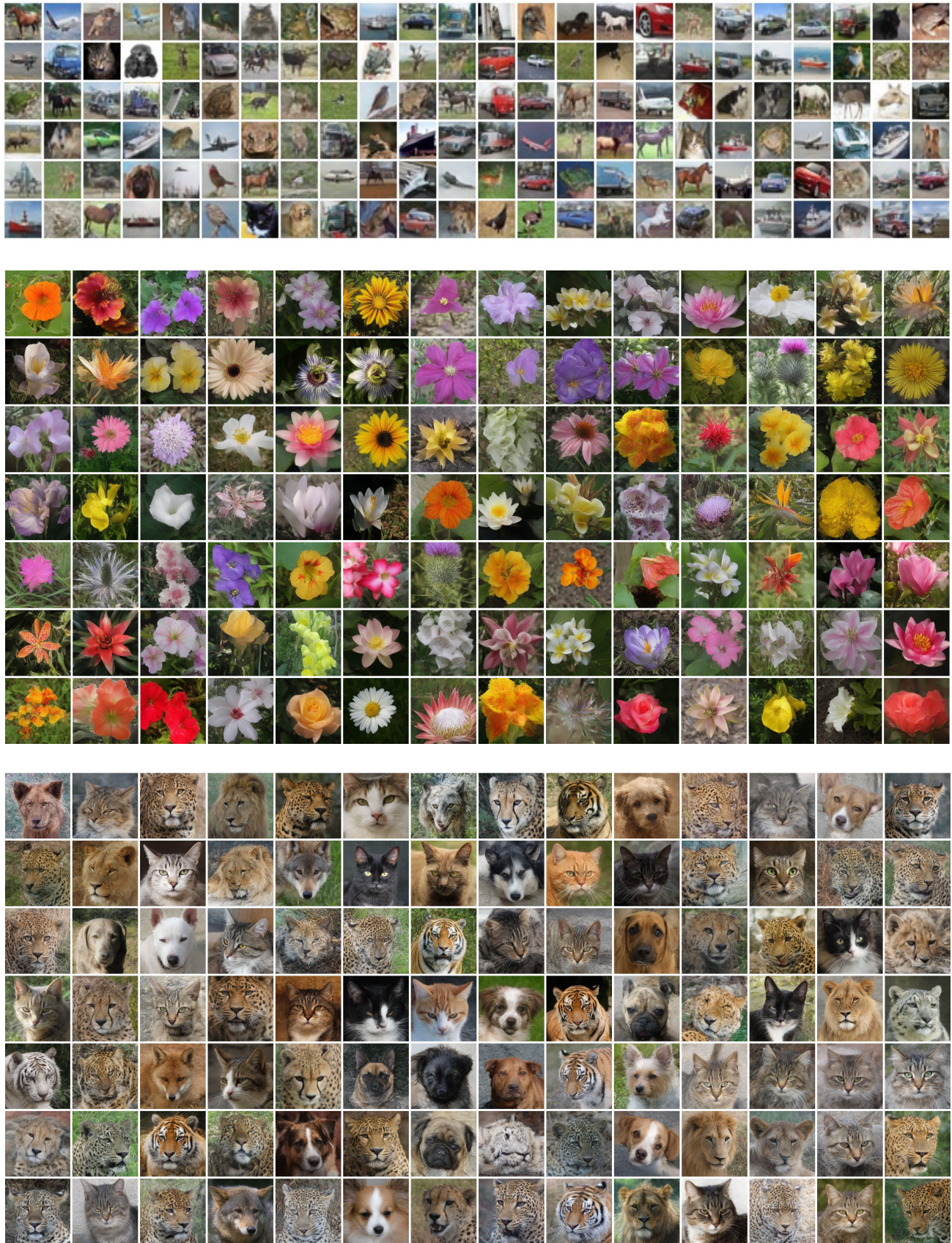


Figure 17. **Uncurated qualitative results** on CIFAR-10 (32×32), Oxford-Flowers and Animal-Faces (256×256). Results are 2-step generation without CFG.

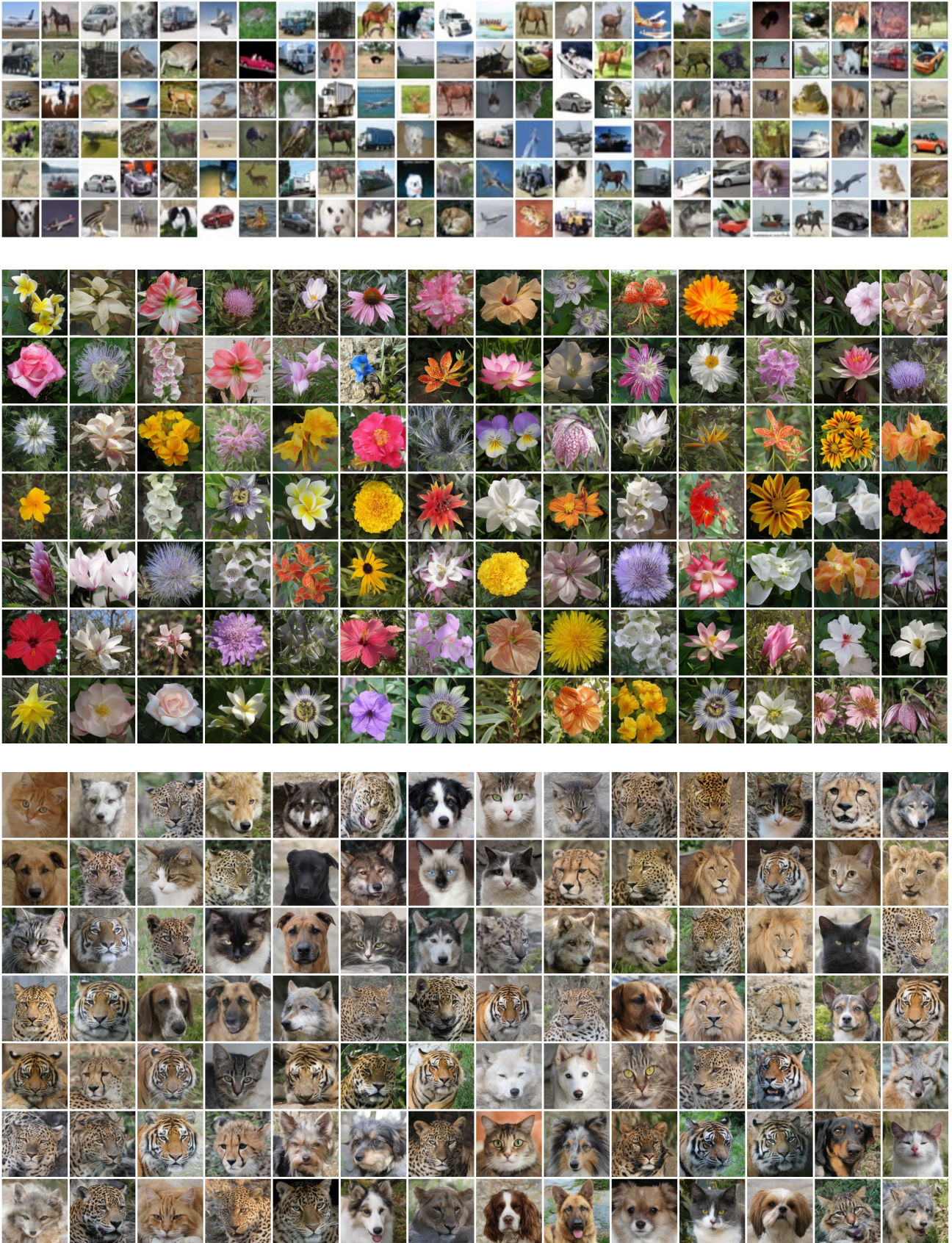


Figure 18. **Uncurated qualitative results** on CIFAR-10 (32×32), Oxford-Flowers and Animal-Faces (256×256). Results are 4-step generation without CFG.

A. Additional Results on CIFAR-10

Table 8 shows that using $\text{CFG} = 1.2$ yields a slight improvement over the baseline in Table 1. For example, with 6-step sampling, gFID decreases from 1.65 (no CFG) to 1.41 (with CFG), and IS increases from 10.7 to 10.8.

B. Memorization Risk on CIFAR-10

In the case of training longer on CIFAR-10, *e.g.*, around 10K epochs following common practice of diffusion models (Song et al., 2023; Geng et al., 2025), we found our model sometimes generates near-duplicate samples. Figure 19 presents near-duplicate samples, *i.e.*, the bird, from different sampling runs. This phenomenon indicates our model may memorize some training samples when trained excessively to overfit the small-scale data distribution. However, our model generates the flipped version of the bird, which is not exactly the same as the real image, suggesting that the model does not simply copy the training image but rather learns a transformation of it. This aligns with known memorization issue in diffusion models (Somepalli et al., 2023; Wen et al., 2024). This longer training improves generation quality further, as shown in Table 7, *e.g.*, gFID = 0.94 and IS = 11.1 with 10-step sampling with CFG.

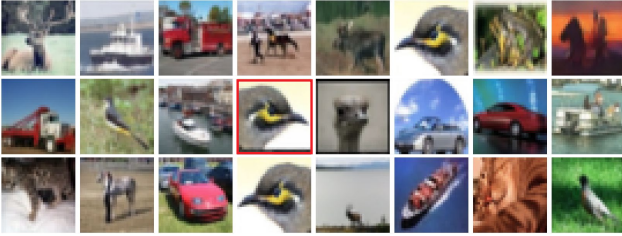


Figure 19. **Memorization risk** with longer 10K training epochs on CIFAR-10. Each row is a different sampling run showing near-duplicate birds of the real training image (in red box).

C. Additional Ablations

C.1. CFG Position

In Algorithm 1, we have three options to apply CFG: (1) only in the pixel space after decoding; (2) only in the latent space after encoding; (3) in both latent and pixel spaces, termed as “combo”. Since the combo option applies CFG twice, we halve the CFG scale for each application to keep the overall strength consistent. For example, if the overall CFG scale is s , we use $s^{1/2}$ for each position. Table 9 presents the results on ImageNet with different CFG positions and scales. We observe that applying CFG in the pixel space consistently outperforms applying it in the latent space. The combo option with $\text{CFG} = 1.6$ achieves the best results for 4-step sampling on ImageNet. Unless otherwise specified, we apply combo for all experiments.

Table 7. **Few-step conditional generation results** on CIFAR-10 with longer 10K training epochs.

METHOD	STEPS	RFD ↓	GFID ↓	IS ↑
WITHOUT CFG				
SPHERE-L	1	0.68	16.02	9.5
	2	-	4.92	10.3
	4	-	1.24	10.7
	6	-	1.01	10.9
	8	-	1.01	10.8
	10	-	1.00	10.9
WITH CFG = 1.2				
SPHERE-L	1	-	16.57	9.5
	2	-	5.17	10.4
	4	-	1.20	10.8
	6	-	0.96	10.9
	8	-	0.94	10.9
	10	-	0.94	11.1

Table 8. **Few-step generation results** on CIFAR-10 with CFG.

METHOD	CFG	STEPS	RFD ↓	GFID ↓	IS ↑
SPHERE-L	1.2	1	0.59	18.82	9.1
		2	-	8.41	10.0
		4	-	2.53	10.6
		6	-	1.41	10.8
		8	-	1.18	10.8
		10	-	1.12	10.9

C.2. Dialing in the Noise Distribution

Because the noise magnitude has a strong impact on image quality, we dial this in by considering a few more sophisticated sampling strategies for the noise magnitude. After determining the optimal maximum angle α in the previous section, we further explore whether mixing a range of larger angles during training can enhance generation quality. We conduct experiments on CIFAR-10 with latent size $L = 16 \times 16 \times 8$. We train the model for 2000 epochs.

In Table 10, we compare three settings: (1) first row: a base case where angles are sampled uniformly from $[0^\circ, 80^\circ]$; (2) second row: the base case mixing with larger angles from $[80^\circ, 85^\circ]$ with 0.1 probability; and (3) third row: the base case mixing with even larger angles from $[80^\circ, 89^\circ]$ with 0.1 probability.

We observed that mixing in a constrained range of larger angles (*e.g.*, $[80^\circ, 85^\circ]$) improves generation quality for both one-step and four-step sampling. However, including excessively high angles, *e.g.*, $[80^\circ, 89^\circ]$, degrades generation quality and also causes unstable training with gradient explosions. We therefore conclude that augmenting the training data with a moderate band of large angles enhances quality.

Accordingly, we adopt this mixing strategy for all experiments: for CIFAR-10 with small image size 32, we add angles from $[80^\circ, 85^\circ]$ with 0.1 probability to the base range of $[0^\circ, 80^\circ]$; for ImageNet, Animal-Faces, Oxford-Flowers with large image size 256, we add angles from $[85^\circ, 89^\circ]$ with 0.1 probability to the base range of $[0^\circ, 85^\circ]$.

Table 9. **Ablation on CFG position.** Results are reported with gFID \downarrow on ImageNet. The sampling scheme is the fixed strength with sharing noise ϵ in spherifying process across all steps.

CFG POSITION	STEPS W/ CFG = 1.0					STEPS W/ CFG = 1.2					STEPS W/ CFG = 1.4					STEPS W/ CFG = 1.6				
	1	2	4	6	8	1	2	4	6	8	1	2	4	6	8	1	2	4	6	8
ENC	16.7	7.9	6.0	7.9	10.0	16.7	8.0	6.1	7.8	10.3	16.8	7.9	6.1	7.9	10.5	16.8	7.9	6.11	8.1	10.7
DEC	16.7	8.0	5.9	7.8	10.1	15.0	8.3	4.7	5.1	6.2	15.4	9.3	5.0	4.5	5.0	16.4	10.2	5.7	4.7	4.8
COMBO	16.7	7.9	6.0	7.8	10.2	15.3	8.0	5.1	6.1	7.9	15.1	8.2	4.7	5.2	6.5	15.0	8.6	4.7	4.7	5.7

Table 10. **Impact of mixing a range of larger angles** during training on CIFAR-10. All settings yield rFID scores in the range of 0.46–0.48.

BASE MAX α	MIX RANGE	STEPS	W/O CFG		W/ CFG = 1.6	
			gFID \downarrow	IS \uparrow	gFID \downarrow	IS \uparrow
80°	-	1	21.17	8.8	18.93	9.3
		4	8.48	10.0	7.47	10.1
80°	[80°, 85°]	1	19.82	8.7	18.39	9.1
		4	7.98	10.1	7.05	10.2
80°	[80°, 89°]	1	29.60	8.8	30.23	9.2
		4	9.38	10.1	8.21	10.1

C.3. Explicit Distribution Regularization

The ideal distribution on the sphere is uniform. To form such a distribution, we could add a regularization term to the encoder output. We investigate two options: (1) Batch Normalization (BN) (Ioffe, 2015) on the encoder output before spherifying; Since we sample noise from a Normal distribution, BN encourages the encoder output to be Gaussian, which is close to the distribution of noise. (2) SWD loss (Rowland et al., 2019; Kolouri et al., 2019; Wu et al., 2019; Deshpande et al., 2019), which is a sliced Wasserstein distance between the encoder output and uniform distribution on the sphere. We implement the loss \mathcal{L}_{SWD} by constructing orthogonal random projections following (Rowland et al., 2019) Algorithm 3 to the encoder output before spherifying.

Starting with the baseline model without any regularization, we then gradually add BN and SWD loss to evaluate their effects on generation quality. As shown in Table 11, applying BN marginally improves generation quality for both one-step and four-step sampling. However, adding SWD loss on top of BN slightly degrades generation quality.

This suggests that our noisy spherifying method already encourages a near-uniform distribution on the sphere, and additional BN or SWD regularization may not be necessary. There are downsides to these regularizations as well. BN introduces extra calibration steps during inference to adjust batch statistics, which complicates the generation process (Brock et al., 2018) (see details in Appendix C.4). And SWD is expensive with large latent dimensions, requiring latent caching and memory for random projection matrices.

Table 11. **Ablation on explicit uniform regularization** on the spherical latent space.

ADD-ONS	STEPS	rFID \downarrow	W/O CFG		W/ CFG = 1.6	
			gFID \downarrow	IS \uparrow	gFID \downarrow	IS \uparrow
-	1	1.20	17.56	182.5	14.88	217.8
	4	-	8.73	177.6	8.05	245.0
+ BN	1	0.94	17.28	181.6	14.71	216.0
	4	-	8.39	171.3	7.65	240.3
+ \mathcal{L}_{SWD}	1	1.01	18.95	164.4	15.47	202.1
	4	-	8.41	165.5	7.46	235.2

Table 12. **Recalibrating BN stats** for inference on ImageNet.

CALIB BY	STEPS	rFID \downarrow	W/O CFG		W/ CFG = 1.6	
			gFID \downarrow	IS \uparrow	gFID \downarrow	IS \uparrow
ON-THE-FLY	1	0.99	14.21	214.2	13.02	246.0
	4	-	6.85	202.3	7.39	268.2
GENERATED IMAGES	1	0.99	14.21	214.2	13.02	246.0
	4	-	7.30	211.1	8.08	273.9
REFERENCE IMAGES	1	1.00	14.21	214.2	13.02	246.0
	4	-	7.46	209.3	8.06	273.6

C.4. BatchNorm Recalibration

A common issue from BatchNorm (Ioffe, 2015) is the train-test discrepancy for the batch statistics, *i.e.*, running mean and variance. Directly using the training statistics during inference may lead to performance degradation. Using batch norm, we need to calibrate the stats of bn layers during inference. DCGANs (Radford et al., 2016) does not calibrate it. However, BigGANs (Brock et al., 2018) found the calibration affect the generation quality significantly, and run the testing samples to get the BN stats.

In Table 12, we compare the results of using training stats and calibrated stats for BN on ImageNet. For calibration, we have three options: (1) no calibration: reset the BN training stats and then directly use the model for generation, the BN stats are updated on-the-fly. (2) calibrate with reference images: run the model on the training set to accumulate BN statistics prior to generation. (3) calibrate with generated images: sample 500 images to get the BN stats before generating the final samples for evaluation. The results indicate that explicit calibration is not strictly necessary.

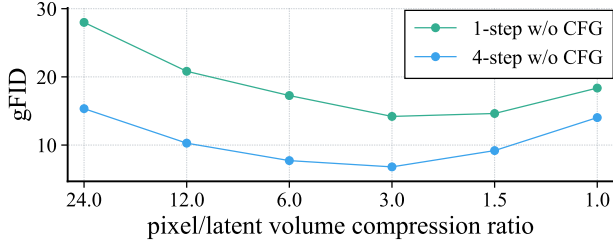


Figure 20. **Quantitative impact of volume compression ratio** on ImageNet. Details in Table 13.

C.5. Volume Compression Ratio

Distinct from classical VAEs, our encoder outputs a high latent dimension. Our pixel/latent volume compression ratio is much smaller than that of a standard VAE, for example, 1.5 for CIFAR-10 and 3.0 for ImageNet, compared to 48 for a standard VAE (*e.g.*, from pixel $256^2 \times 3$ to latent $32^2 \times 4$) (Rombach et al., 2022; HaCohen et al., 2024). We further study the impact of pixel/latent volume compression ratio for the latent dim L . It is well known that, for autoencoders, a higher compression ratio typically leads to worse reconstruction quality but eases diffusion models to fit in (Yao et al., 2025; Zheng et al., 2025). However, this is not necessarily true for our method. We vary the channel depth of latent dim L to adjust the volume compression ratio while keeping the spatial resolution fixed, *i.e.*, 16^2 for image size 256 on ImageNet. The results are plotted in Figure 20, demonstrating two optimal compression ratios around 1.5 and 3.0 for few-step generation. This ratio is way lower than that of typical autoencoders used in diffusion models, *e.g.*, 48 in (Rombach et al., 2022; Zheng et al., 2025). Unless otherwise specified, we use the optimal compression ratio 3.0 for ImageNet and 1.5 for CIFAR-10, Animal-Faces and Oxford-Flowers.

C.6. Noise Prior Distribution

Algorithm 1 samples noise \mathbf{e} in Equation (3) from a isotropic Gaussian $\mathcal{N}(0, I)$ as the “input latent” for generation. Since we spherify it before feeding into the decoder, we assume the generation should be insensitive to the specific choice of noise prior distribution. This is contrast to GANs (Brock et al., 2019; Karras et al., 2020; Sauer et al., 2022), which apply truncation tricks (Karras et al., 2019) to the latent prior distribution to improve generation fidelity. We compare two noise prior distributions: (1) the standard Gaussian $\mathcal{N}(0, I)$; and (2) a truncated Normal $\mathcal{N}(0, I)$ truncated to $[-\alpha, \alpha]$, where α is the truncation threshold. Table 14 confirms the generation quality is insensitive to the choice of noise prior distribution, as the gFID and IS are similar across different α values.

Table 13. **Ablation on pixel/latent volume compression ratio** on ImageNet without CFG.

COMPRESSION RATIO	LATENT L	STEPS	RFID ↓	gFID ↓	IS ↑
24	$16^2 \times 32$	1	1.61	27.99	195.6
		4	-	15.34	245.1
12	$16^2 \times 64$	1	1.38	20.82	233.8
		4	-	10.28	259.0
6.0	$16^2 \times 128$	1	1.15	17.26	225.0
		4	-	7.72	229.4
3.0	$16^2 \times 256$	1	1.01	14.20	214.2
		4	-	6.80	201.7
1.5	$16^2 \times 512$	1	0.69	14.63	136.0
		4	-	9.19	125.9
1.0	$16^2 \times 768$	1	0.64	18.36	97.6
		4	-	14.03	85.6

Table 14. **Ablation of noise prior distribution** on ImageNet.

TRUNC α	STEPS	W/O CFG		W/ CFG = 1.2	
		gFID ↓	IS ↑	gFID ↓	IS ↑
-	1	17.81	233.9	16.30	261.0
	4	5.92	192.3	5.02	229.7
0.05	1	17.89	231.7	15.81	277.0
	4	5.88	192.2	5.08	228.2
0.5	1	17.96	231.4	15.73	275.8
	4	5.89	193.5	5.02	229.7
2.0	1	17.82	233.5	15.66	281.0
	4	5.85	191.7	5.06	228.3

Table 15. **Detailed quantitative results for the impact of the angle α** on ImageNet for Figure 10.

ANGLE α	σ	STEPS	RFID ↓	W/O CFG		W/ CFG = 1.6	
				gFID ↓	IS ↑	gFID ↓	IS ↑
45°	1.0	1	0.69	69.57	31.3	59.09	39.9
		4	-	80.32	18.7	73.04	21.9
80°	5.7	1	0.87	60.14	33.4	51.73	41.3
		4	-	42.72	39.3	34.14	54.3
81°	6.3	1	0.94	47.68	40.4	40.11	51.1
		4	-	37.64	43.4	28.88	62.2
83°	8.1	1	0.99	20.15	107.7	16.21	135.1
		4	-	13.05	102.2	8.70	151.0
84°	9.5	1	1.15	13.96	161.7	12.30	186.0
		4	-	11.92	116.9	8.66	171.5
85°	11.4	1	1.32	15.40	188.3	13.14	225.1
		4	-	7.53	176.4	7.28	243.1
86°	14.3	1	1.26	15.37	232.7	14.15	267.1
		4	-	9.96	229.7	11.14	292.0
87°	19.1	1	1.29	16.40	245.0	14.94	291.3
		4	-	9.99	257.2	11.50	315.1
88°	28.6	1	1.53	23.27	203.5	17.72	268.7
		4	-	13.15	241.5	12.24	304.4

Table 16. Training hyperparameters and model details for main experiments.

HYPERPARAMETER \ DATASET	CIFAR-10	ANIMAL-FACES	OXFORD-FLOWERS	IMAGENET
IMAGE SIZE	32	256	256	256
BATCH SIZE	256	256	256	256
LEARNING RATE (LR)	1×10^{-4}	1×10^{-4}	1×10^{-4}	1×10^{-4}
LR DECAY SCHEDULE	COSINE	COSINE	COSINE	COSINE
MIN LR	1×10^{-6}	1×10^{-6}	1×10^{-6}	1×10^{-6}
WEIGHT DECAY	0.0	0.0	0.0	0.0
OPTIMIZER	ADAMW	ADAMW	ADAMW	ADAMW
WARMUP EPOCHS	10	10	10	5
TOTAL EPOCHS	5K	1K	1K	800
MODEL	SPHERE-L	SPHERE-L	SPHERE-L	SPHERE-L, -XL
ENCODER / DECODER SIZE	LARGE	LARGE	LARGE	LARGE, XLARGE
NUMBER OF TRANSFORMER BLOCKS	24	24	24	24, 28
NUMBER OF ATTENTION HEADS	16	16	16	16
TRANSFORMER HIDDEN SIZE	1024	1024	1024	1024, 1152
MLP-MIXER DEPTH	2	4	4	4
CLASS CONDITION	✓	-	✓	✓
MODEL PARAMS	921M	642M	948M	950M
VOLUME COMPRESSION RATIO	1.5	1.5	1.5	3.0
LATENT DIM L	$16^2 \times 8$	$32^2 \times 128$	$32^2 \times 128$	$32^2 \times 64$
ANGLE α JITTER RANGE	$[0^\circ, 80^\circ]$	$[0^\circ, 85^\circ]$	$[0^\circ, 85^\circ]$	$[0^\circ, 85^\circ]$
ANGLE α MIX RANGE	$[80^\circ, 85^\circ]$	$[85^\circ, 89^\circ]$	$[85^\circ, 89^\circ]$	$[85^\circ, 89^\circ]$
ANGLE α MIX PROBABILITY	0.1	0.1	0.1	0.1
$\mathcal{L}_{\text{PIX-RECON}}$ SMOOTH L1 LOSS WEIGHT	1.0	25.0	25.0	50.0
$\mathcal{L}_{\text{PIX-RECON}}$ PERCEPTUAL LOSS WEIGHT	1.0	1.0	1.0	1.0
$\mathcal{L}_{\text{PIX-CON}}$ SMOOTH L1 LOSS WEIGHT	0.5	1.0	1.0	25.0
$\mathcal{L}_{\text{PIX-CON}}$ PERCEPTUAL LOSS WEIGHT	0.5	1.0	1.0	1.0
$\mathcal{L}_{\text{LAT-CON}}$ LOSS WEIGHT	0.1	0.1	0.1	0.1

D. Hyperparameters

Table 16 lists the training hyperparameters and model details for our main experiments. In addition, for unconditional generation on CIFAR-10, the only difference is removing the class condition and training for 10K epochs. We found EMA for smoothing weights in checkpoints may not noticeably change FID or sample quality, likely because we use cosine annealing learning rate schedule, which already provides a smoothing effect.

Impact Statement

Our work proposes a new generative framework that offers a fresh perspective on image generation and yields various benefits, including fast sampling and high-dimensional latent space for image generation. Although this specific method does not raise unique ethical challenges, we acknowledge ongoing general concerns inherent to the field, and we encourage continued community discussion to ensure responsible development and mitigation of potential risks.

References

- Aneja, J., Schwing, A., Kautz, J., and Vahdat, A. A contrastive learning approach for training variational autoencoder priors. In *NeurIPS*, 2021.
- Bourlard, H. and Kamp, Y. Auto-association by multilayer perceptrons and singular value decomposition. *Biological cybernetics*, 1988.
- Brock, A., Donahue, J., and Simonyan, K. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2018.
- Brock, A., Donahue, J., and Simonyan, K. Large scale gan training for high fidelity natural image synthesis. In *ICLR*, 2019.
- Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., and Sutskever, I. Generative pretraining from pixels. In *ICML*, 2020.
- Child, R. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Choi, Y., Uh, Y., Yoo, J., and Ha, J.-W. Stargan v2: Diverse image synthesis for multiple domains. In *CVPR*, 2020.
- Dai, B. and Wipf, D. Diagnosing and enhancing vae models. In *ICLR*, 2019.
- Davidson, T. R., Falorsi, L., De Cao, N., Kipf, T., and Tomczak, J. M. Hyperspherical variational auto-encoders. *arXiv preprint arXiv:1804.00891*, 2018.
- De Cao, N. and Aziz, W. The power spherical distribution. *arXiv preprint arXiv:2006.04437*, 2020.
- Denton, E. L., Chintala, S., Fergus, R., et al. Deep generative image models using a Laplacian pyramid of adversarial networks. In *NeurIPS*, 2015.
- Deshpande, I., Hu, Y.-T., Sun, R., Pyrros, A., Siddiqui, N., Koyejo, S., Zhao, Z., Forsyth, D., and Schwing, A. G. Max-sliced wasserstein distance and its use for gans. In *CVPR*, 2019.
- Dhariwal, P. and Nichol, A. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021.
- Dosovitskiy, A. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020.
- Esser, P., Kulal, S., Blattmann, A., Entezari, R., Müller, J., Saini, H., Levi, Y., Lorenz, D., Sauer, A., Boesel, F., et al. Scaling rectified flow transformers for high-resolution image synthesis. In *ICML*, 2024.
- Frans, K., Hafner, D., Levine, S., and Abbeel, P. One step diffusion via shortcut models. *arXiv preprint arXiv:2410.12557*, 2024.
- Geng, Z., Pockle, A., Luo, W., Lin, J., and Kolter, J. Z. Consistency models made easy. *arXiv preprint arXiv:2406.14548*, 2024.
- Geng, Z., Deng, M., Bai, X., Kolter, J. Z., and He, K. Mean flows for one-step generative modeling. In *NeurIPS*, 2025.
- Ghosh, P., Sajjadi, M. S., Vergari, A., Black, M., and Schölkopf, B. From variational to deterministic autoencoders. In *ICLR*, 2020.
- Girshick, R. Fast r-cnn. In *ICCV*, 2015.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *NeurIPS*, 2014.
- HaCohen, Y., Chiprut, N., Brazowski, B., Shalem, D., Moshe, D., Richardson, E., Levin, E., Shiran, G., Zabari, N., Gordon, O., et al. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024.
- Hawthorne, C., Jaegle, A., Cangea, C., Borgeaud, S., Nash, C., Malinowski, M., Dieleman, S., Vinyals, O., Botvinick, M., Simon, I., et al. General-purpose, long-context autoregressive modeling with percever ar. In *ICML*, 2022.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, volume 30, 2017.
- Hinton, G. E. and Zemel, R. Autoencoders, minimum description length and helmholtz free energy. In *NeurIPS*, 1993.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2022.

- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- Hoogetboom, E., Heek, J., and Salimans, T. simple diffusion: End-to-end diffusion for high resolution images. In *ICML*, 2023.
- Hoogetboom, E., Mensink, T., Heek, J., Lamerigts, K., Gao, R., and Salimans, T. Simpler diffusion (sid2): 1.5 fid on imagenet512 with pixel-space diffusion. *arXiv preprint arXiv:2410.19324*, 2024.
- Ioffe, S. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- Jayasumana, S., Ramalingam, S., Veit, A., Glasner, D., Chakrabarti, A., and Kumar, S. Rethinking fid: Towards a better evaluation metric for image generation. In *CVPR*, 2024.
- Johnson, J., Alahi, A., and Fei-Fei, L. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016.
- Kang, M., Zhu, J.-Y., Zhang, R., Park, J., Shechtman, E., Paris, S., and Park, T. Scaling up gans for text-to-image synthesis. In *CVPR*, 2023.
- Karras, T., Aila, T., Laine, S., and Lehtinen, J. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*, 2018.
- Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019.
- Karras, T., Aittala, M., Hellsten, J., Laine, S., Lehtinen, J., and Aila, T. Training generative adversarial networks with limited data. In *NeurIPS*, 2020.
- Ke, G. and Xue, H. Hyperspherical latents improve continuous-token autoregressive generation. *arXiv preprint arXiv:2509.24335*, 2025.
- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. In *NeurIPS*, 2018.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Kingma, D. P. and Welling, M. An introduction to variational autoencoders. *arXiv preprint arXiv:1906.02691*, 2019.
- Kolouri, S., Pope, P. E., Martin, C. E., and Rohde, G. K. Sliced-wasserstein autoencoder: An embarrassingly simple generative model. In *ICLR*, 2019.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. *Toronto, ON, Canada*, 2009.
- Kumar, A. and Patel, V. M. Learning on the manifold: unlocking standard diffusion transformers with representation encoders. *arXiv preprint arXiv:2602.10099*, 2026.
- Labs, B. F., Batifol, S., Blattmann, A., Boesel, F., Consul, S., Diagne, C., Dockhorn, T., English, J., English, Z., Esser, P., et al. Flux. 1 kontext: Flow matching for in-context image generation and editing in latent space. *arXiv preprint arXiv:2506.15742*, 2025.
- LeCun, Y. Phd thesis: Modeles connexionnistes de l'apprentissage (connectionist learning models). <https://api.semanticscholar.org/CorpusID:151887454>, 1987.
- Li, T. and He, K. Back to basics: Let denoising generative models denoise. *arXiv preprint arXiv:2511.13720*, 2025.
- Li, T., Sun, Q., Fan, L., and He, K. Fractal generative models. In *TMLR*, 2025.
- Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and Le, M. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Liu, X., Gong, C., and Liu, Q. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *ICLR*, 2023.
- Lu, C. and Song, Y. Simplifying, stabilizing and scaling continuous-time consistency models. In *ICLR*, 2025.
- Ma, N., Goldstein, M., Albergo, M. S., Boffi, N. M., Vanden-Eijnden, E., and Xie, S. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In *ECCV*, 2024.
- Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I., and Frey, B. Adversarial autoencoders. In *ICLR*, 2015.
- Nichol, A. Q., Dhariwal, P., and et al. Improved denoising diffusion probabilistic models. In *ICML*, 2021.
- Nilsback, M.-E. and Zisserman, A. Automated flower classification over a large number of classes. In *ICCV*, 2008.
- Peebles, W. and Xie, S. Scalable diffusion models with transformers. In *ICCV*, 2023.
- Perez, E., Strub, F., De Vries, H., Dumoulin, V., and Courville, A. Film: Visual reasoning with a general conditioning layer. In *AAAI*, 2018.
- Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. Sdxl: Improving latent diffusion models for high-resolution image synthesis. In *ICLR*, 2024.

- Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Rezende, D. J. and Viola, F. Taming vaes. *arXiv preprint arXiv:1810.00597*, 2018.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- Rowland, M., Hron, J., Tang, Y., Choromanski, K., Sarlos, T., and Weller, A. Orthogonal estimation of wasserstein distances. In *AISTATS*, 2019.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. In *IJCV*, 2015.
- Salimans, T. and Ho, J. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. Improved techniques for training gans. In *NeurIPS*, 2016.
- Sauer, A., Schwarz, K., and Geiger, A. Stylegan-xl: Scaling stylegan to large diverse datasets. In *SIGGRAPH*, 2022.
- Sauer, A., Lorenz, D., Blattmann, A., and Rombach, R. Adversarial diffusion distillation. In *ECCV*, 2024.
- Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and Goldstein, T. Understanding and mitigating copying in diffusion models. In *NeurIPS*, 2023.
- Song, J., Meng, C., and Ermon, S. Denoising diffusion implicit models. In *ICLR*, 2021.
- Song, Y., Dhariwal, P., Chen, M., and Sutskever, I. Consistency models. In *ICML*, 2023.
- Stein, G., Cresswell, J., Hosseinzadeh, R., Sui, Y., Ross, B., Vilecroze, V., Liu, Z., Caterini, A. L., Taylor, E., and Loaiza-Ganem, G. Exposing flaws of generative model evaluation metrics and their unfair treatment of diffusion models. In *NeurIPS*, 2023.
- Studer, C. and Bölcskei, H. Soft-input soft-output single tree-search sphere decoding. *IEEE Transactions on Information Theory*, 2010.
- Studer, C., Burg, A., and Bölcskei, H. Soft-output sphere decoding: Algorithms and vlsi implementation. *IEEE Journal on selected areas in Communications*, 2008.
- Su, J., Ahmed, M., Lu, Y., Pan, S., Bo, W., and Liu, Y. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 2024.
- Tian, K., Jiang, Y., Yuan, Z., Peng, B., and Wang, L. Visual autoregressive modeling: Scalable image generation via next-scale prediction. In *NeurIPS*, 2024.
- Tolstikhin, I., Bousquet, O., Gelly, S., and Schoelkopf, B. Wasserstein auto-encoders. In *ICLR*, 2018.
- Tolstikhin, I. O., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., Lucic, M., and Dosovitskiy, A. Mlp-mixer: An all-mlp architecture for vision. In *NeurIPS*, 2021.
- Tomczak, J. and Welling, M. Vae with a vampprior. In *AISTATS*, 2018.
- Tong, S., Zheng, B., Wang, Z., Tang, B., Ma, N., Brown, E., Yang, J., Fergus, R., LeCun, Y., and Xie, S. Scaling text-to-image diffusion transformers with representation autoencoders. *arXiv preprint arXiv:2601.16208*, 2026.
- Tschannen, M., Pinto, A. S., and Kolesnikov, A. Jetformer: An autoregressive generative model of raw images and text. In *ICLR*, 2025.
- Van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al. Conditional image generation with pixelcnn decoders. In *NeurIPS*, 2016.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *NeurIPS*, 2017.
- Wan, T., Wang, A., Ai, B., Wen, B., Mao, C., Xie, C.-W., Chen, D., Yu, F., Zhao, H., Yang, J., et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- Wen, Y., Liu, Y., Chen, C., and Lyu, L. Detecting, explaining, and mitigating memorization in diffusion models. In *ICLR*, 2024.
- Wu, J., Huang, Z., Acharya, D., Li, W., Thoma, J., Paudel, D. P., and Gool, L. V. Sliced wasserstein generative models. In *CVPR*, 2019.
- Xie, S., Xiao, Z., Kingma, D., Hou, T., Wu, Y. N., Murphy, K. P., Salimans, T., Poole, B., and Gao, R. Em distillation for one-step diffusion models. In *NeurIPS*, 2024.

- Xu, J. and Durrett, G. Spherical latent spaces for stable variational autoencoders. In *EMNLP*, 2018.
- Yang, L., Zhang, Z., Zhang, Z., Liu, X., Xu, M., Zhang, W., Meng, C., Ermon, S., and Cui, B. Consistency flow matching: Defining straight flows with velocity consistency. *arXiv preprint arXiv:2407.02398*, 2024.
- Yao, J., Yang, B., and Wang, X. Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models. In *CVPR*, 2025.
- Yin, T., Gharbi, M., Zhang, R., Shechtman, E., Durand, F., Freeman, W. T., and Park, T. One-step diffusion with distribution matching distillation. In *CVPR*, 2024.
- Yu, L., Simig, D., Flaherty, C., Aghajanyan, A., Zettlemoyer, L., and Lewis, M. Megabyte: Predicting million-byte sequences with multiscale transformers. In *NeurIPS*, 2023.
- Yu, S., Kwak, S., Jang, H., Jeong, J., Huang, J., Shin, J., and Xie, S. Representation alignment for generation: Training diffusion transformers is easier than you think. *arXiv preprint arXiv:2410.06940*, 2024.
- Yu, Y., Xiong, W., Nie, W., Sheng, Y., Liu, S., and Luo, J. Pixeldit: Pixel diffusion transformers for image generation. *arXiv preprint arXiv:2511.20645*, 2025.
- Zhang, B. and Sennrich, R. Root mean square layer normalization. In *NeurIPS*, 2019.
- Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., and Metaxas, D. N. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017.
- Zhang, R., Isola, P., Efros, A. A., Shechtman, E., and Wang, O. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- Zhao, D., Zhu, J., and Zhang, B. Latent variables on spheres for autoencoders in high dimensions. *arXiv preprint arXiv:1912.10233*, 2019.
- Zheng, B., Ma, N., Tong, S., and Xie, S. Diffusion transformers with representation autoencoders. *arXiv preprint arXiv:2510.11690*, 2025.