# Generative Adversarial Networks via a Composite Annealing of Noise and Diffusion

Kensuke Nakamura

*Computer Science Department, Chung-Ang University, Seoul, Korea*

Simon Korman

*Department of Computer Science, University of Haifa, Israel*

Byung-Woo Hong*

*Computer Science Department, Chung-Ang University, Seoul, Korea*

**Abstract**

Generative adversarial network (GAN) is a framework for generating fake data using a set of real examples. However, GAN is unstable in the training stage. In order to stabilize GANs, the noise injection has been used to enlarge the overlap of the real and fake distributions at the cost of increasing variance. The diffusion (or smoothing) may reduce the intrinsic underlying dimensionality of data but it suppresses the capability of GANs to learn high-frequency information in the training procedure. Based on these observations, we propose a data representation for the GAN training, called noisy scale-space (NSS), that recursively applies the smoothing with a balanced noise to data in order to replace the high-frequency information by random data, leading to a coarse-to-fine training of GANs. We experiment with NSS using DCGAN and StyleGAN2 based on benchmark datasets in which the NSS-based GANs outperforms the state-of-the-arts in most cases.

*Keywords:* generative adversarial networks, optimization, scale-space, noise injection, coarse-to-fine training

*Corresponding author
Email address:* `hong@cau.ac.kr` (Byung-Woo Hong)

## 1. Introduction

Generative adversarial network (GAN) [1] is a machine learning framework to generate realistic fake data. GAN learns the probabilistic distribution of the training (real) data using two adversarial networks: the generator that is trained to create realistic fake data from a random seed called the latent vector, and the discriminator that is dedicated to distinguish the reals against fake data. GAN has been studied extensively in the several past years and currently is an essential tool for a wide variety of applications.

However, the training procedure is prone to numerical unstability in GANs [2, 3, 4, 5, 6]. Since it is a two-player game between the generator and discriminator [4], the optimization can fall into a local minima where the discriminator reaches a perfect solution first and the generator cannot be trained anymore. This failure of GAN severely limits the quality of fake data, resulting in the mode collapse. The failure of GAN occurs when there is almost no overlap between the real distribution with fake distribution [7], in particular, in the beginning of training when the discriminator can reject fake data with a high confidence [1]. Therefore, GANs require stabilization methods in the optimization process in order to obtain better fake data.

The stabilization methods of GAN are different with those for the feed-forward deep networks, e.g., weight-decay [8] and the momentum [9], due to the dynamics of the discriminator with generator. Since the failure of GAN can arise from the use of KL-divergence in loss, a variety of the discrepancy measures have been studied, e.g., [10, 11, 12], including Wasserstein distance [7, 13] that provides gradients to the generator even with a small overlap of the two distributions. Albeit, the Wasserstein loss is often inferior to the original loss in the quality of fake data [3, 5]. Another strategy is to inhibit discriminator training based on the loss. To this aim, the gradient regularization penalizes gradients of the discriminator. However, it depends on the model that varies during the training [5].

This paper focuses on data-based stabilization of the GAN training that

manipulates only the real and fake data independent of the model architecture. For instance, the repetition of noise injection to both real and fake data (namely noise-space) enlarges the overlap between probability distributions [14]. However, the noise increases the variance of data inevitably. The repetition of data smoothing based on Gaussian kernels, or the scale-space, is a general technique in machine learning that suppresses the high frequency features such as textures and details in the image so as to make each data more simple to learn by algorithms. In the case of GAN, it is expected that the diffusion makes the real data easy to mimic by the generator. However, we have found that it limits the learning capability of the generator to learn high-frequency information.

Based on these observations, we propose an algorithm for stabilizing the optimization of GANs based on a *noisy scale-space* (NSS) that continuously removes high-frequency information in image while adding noise. The proposed noisy scale-space enables a coarse-to-fine training of GANs in which we can train a generative model using low-level information in data with noise without the increment of data variance, while keeping the current model capable of learning the high-level information. We also present a synthetic dataset using the Hadamard bases [15] that can visualize the true distributions of real and fake data in order to characterize the drawback of the conventional scale-space in GAN optimization. Then, we perform experiments with the proposed NSS using DCGAN [16]. The experimental results based on the major datasets show that the proposed NSS-GAN outperforms other methods in most cases irrespective of the image generation tasks. Specifically, it is shown that the stabilization effect by our method is not the simple summation of the noise-space with the scale-space but due to the use of their mutually complementary relationship. We also demonstrate that the proposed NSS can improve the accuracy of StyleGAN2 [17] for high-resolution images.

We relate our method to prior works in Section 2. Then we consider a data-based stabilization of GAN training in Section 3, followed by the proposed noisy scale-space in Section 4. The effectiveness of our method is demonstrated experimentally in Section 5 and we conclude in Section 6.

3

## 2. Related works

**Loss-based GAN stabilization:** The regularization term in the loss has been studied to restrict the update of discriminator. In a related context, Wasserstein-GAN (WGAN) uses a weight clipping [13] to guarantee the Lipschitz constraint. WGAN-gp uses a gradient penalty [18] that can improve the quality of fakes in practice. The spectral normalization [19] is an efficient variant of the gradient-penalty. Dragan [20] penalizes the sharp gradient of discriminator to real data. We will use the non-saturating loss [1] based on the KL-divergence that is known to be the best choice in practice. The non-saturating loss can be further improved using the gradient regularization [3] that penalizes the gradient norm of the discriminator. It is shown that the penalty based on the gradient-norm of the discriminator is equivalent to adding input noise in GAN using $f$-divergence [21]. The drawback of gradient regularization is that it depends on the distribution of fakes determined by the generator which changes during training [5].

**Procedure-based GAN stabilization:** The two time-scale update rule [22] is a method of using different annealings for discriminator and generator in order to slow the convergence of the discriminator. The simultaneous update of the two networks was studied in [23]. Progressive augmentation of GAN [24] extends the label noise [25] into the GAN framework to perturb the real and fake labels. The one-sided label smoothing replaces the 0 and 1 target labels for the discriminator with smoothed values, like 0.9 or 0.1 [18, 26].

**Data-based GAN stabilization:** The proposed method can be categorized in data-based methods that manipulate only data. Lens-GAN [27] introduces a filtering network that transforms the real data against the discriminator, and therefore it still depends on the networks.

The multi-resolution training of GAN is a topic studied in, e.g., Progressive-GAN [28] and MSG-GAN [29], for generating high-resolution images. It trains a shallow network first using low-resolution images, and gradually increases both the number of network layers and the resolution of data. The drawback of the

multi-resolution training is that it strongly limits the architecture.

The data augmentation is another recent topic in GAN training [30, 17, 31] in which a multitude of data transformations, typically consist of spatial transformations (e.g., image rotation, flipping, and cropping) with color transformations (e.g., channel permutation and hue rotation) are combined and applied to both real and fake data in order to enlarge the variation of data. The early works [30, 17] are oriented to avoid over-fitting of GANs to a small set of training examples while the recent works aim to improve the accuracy as we do. Also DistAug [32] has considered a mixture of data transformations in a contrastive training of GANs [33]. However, the mixture of transformations in these studies is a black-box. In contrast, we present a deep understanding of the noise injection and the image diffusion in GAN optimization and then propose a better use of their mutually complementary relationship.

Our method is closely related to the noise injection and the data smoothing. On one hand, the noise injection flattens the probability distribution of data. Thus imposing noise to the real data [34, 7] or both the real and fake data [14] enlarges the overlap of the two distributions. However, adding high-dimensional noise introduces significant variance in the parameter estimation, slowing down the training and requiring multiple samples for counteraction [21, 24]. On the other hand, the scale-space decreases the dimensionality of data by removing high-frequency information and makes data easy to learn by algorithms [35, 36]. However, the conventional scale-space has a critical side-effect in GAN optimization. Different from these baselines, we present a data representation that mitigates the side effects of the noise injection and the data smoothing.

Our algorithm also relates to Ambient-GAN [37] that considers a problem in which incomplete real data containing a Gaussian blurring with additive noise are given. However, Ambient-GAN aims to approximate the original distribution from the incomplete measurements using a conditional network, and created fake images are degenerated by the noise and smoothing. In contrast, we propose a continuous data representation to train GANs, achieving a better quality of fake data than the baseline method using the complete examples.

## 3. Preliminary

### 3.1. Generative adversarial networks

Let us begin with a technical introduction to the generative adversarial networks (GAN). Given a set of real data $(x)$, GAN aims to generate new data with similar statistics as the real data. GAN consists of the two networks: the generator $(G)$ that creates fake data $(G(z))$ from a random latent vector $(z)$, and the discriminator $(D)$ that distinguishes the real data against the fake data. To this aim, a min-max loss $(F)$ is defined [1] as

$$F_{D,G}\big(x, G(z)\big) \quad \coloneqq \quad \mathbb{E}_{x \sim p_{\mathrm{data}}(x)}\big[\log D(x)\big] + \mathbb{E}_{z \sim p_{\mathrm{z}}(z)}\big[\log(1 - D(G(z)))\big], \quad (1)$$

where $p_{\mathrm{data}}(x)$ is the true distribution of the reals, and $p_{\mathrm{z}}(z)$ is a random probability distribution. GANs are optimized using a stochastic gradient descent, e.g., Adam [38], in an alternative way such that the generator is updated to minimize $F$ while the discriminator is trained to maximize $F$.

### 3.2. Data-based stabilization of GAN training

We consider a data-based stabilization method, similar to [30, 17, 31], that projects both the reals and fakes to a data space and feed into the GAN loss as

$$F_{D,G}\big(\Phi_t(x), \Phi_t(G(z))\big), \qquad (2)$$

where $\Phi_t(x)$ and $\Phi_t(G(z))$ are the projected data of reals and fakes using a function $\Phi$ with the time parameter $t$. Using Eq.(2), the discriminator is trained to distinguish the projected real against the projected fake while the generator is trained to create fake data such that the projected fakes $\Phi_t(G(z))$ are similar to the projected reals $\Phi_t(x)$. The data-based stabilization effect will be imposed via a discrete representation of data $\{\Phi_T(y), \Phi_{T-1}(y), ..., \Phi_2(y), \Phi_1(y), \Phi_0(y)\}$ where $\Phi_t(y)$ with larger $t$ is expected to have larger effect such that the data is more easy to learn by GANs, and $\Phi_0(y) = y$ denotes the original real or fake data. To this aim, the time parameter $(t)$ should start with a large value $T$ and shrink to zero during the training, and the choice of $\Phi$ determines the stabilization effect.

(a) NS

(b) SS

(c) NSS

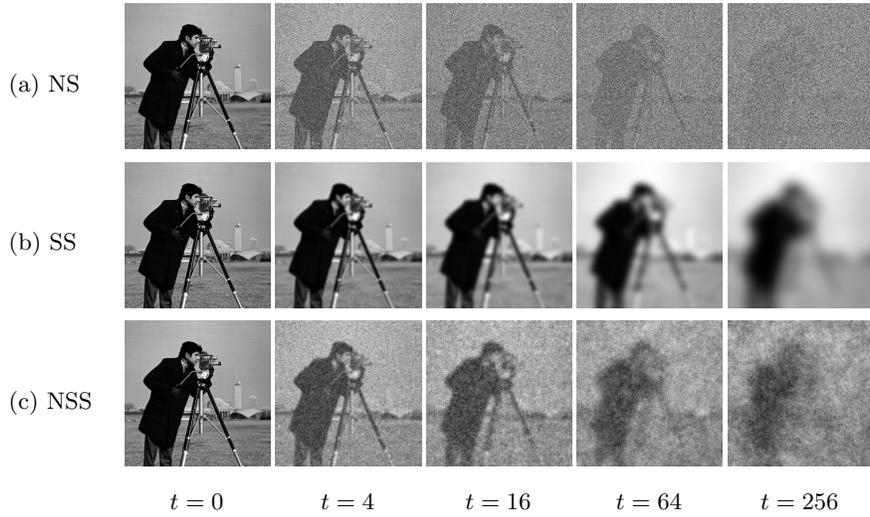$t = 0 \qquad t = 4 \qquad t = 16 \qquad t = 64 \qquad t = 256$

Figure 1: A real image in (a) the conventional noise-space (NS) or the repetition of the Gaussian noise with $\sigma = 0.15$, (b) the conventional scale-space (SS) or the repetition of smoothing, and (c) the proposed noisy scale-space (NSS) (columns) over the time ($t$).
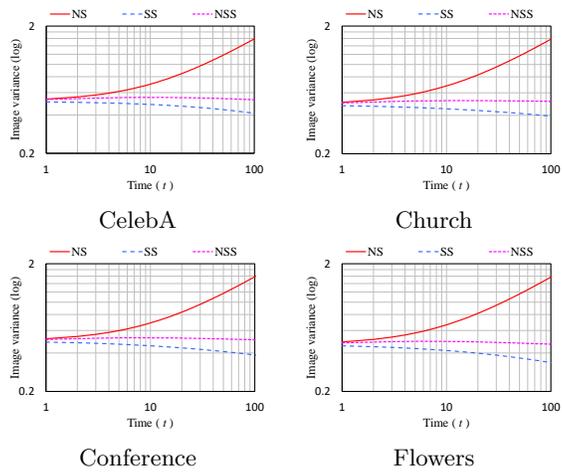


CelebA

Church

Conference

Flowers

Figure 2: (y-axis) The variance of pixel intensities in the noise-space (red line), the scale-space (blue dashed-line), and the proposed noisy scale-space (magenta dotted-line) over (x-axis) the time $t$ based on 128 images in CelebA [39], LSUN-Church, LSUN-Conference [40], and Oxford-Flowers [41] using the noise of $\sigma = 0.15$.

## 4. Stabilization of GAN training via noisy scale-space

*4.1. Proposed noisy scale-space*

We present a data representation, namely noisy scale-space (NSS), that is designed to improve the stability of the optimization for GANs while preserving characteristic features in the generation process. The presented data representation is a balanced composition of the noise with the diffusion such that we can train GANs using the *smoothed* data that is easy to create by the generator first, while flattening the fake distribution in the high-frequency domain, making the current solution be capable of learning the high-frequency information in the further steps. Formally, the noisy scale-space is given as

$$\Phi_t(y) := k * \Phi_{t-1}(y) + \epsilon_t, \tag{3}$$

where $k$ is typically the $3 \times 3$ Gaussian kernel, the symbol $*$ denotes the convolution, $\epsilon_t \sim N(0, \sigma)$, $\Phi_0(y) = y$, and $t \in [T]$, i.e., we apply the kernel and the noise to data simultaneously $T$-times. Figure 1 presents an example of image in (top) the conventional noise-space, (middle) the conventional scale-space, and (bottom) the proposed noisy scale-space. Figure 2 shows the image variance of the three data-spaces.

The noise variance ($\sigma$) is the primal hyper-parameter that balances the diffusion with the noise in the proposed data representation. The conventional scale-space (diffusion) removes high-frequency information in image and decreases the image variance as shown in Figure 2. We determine $\sigma$ in the noisy scale-space such that the image variance is almost preserved over $t$. The noisy scale-space is robust to $\sigma$ for natural images as shown in Figure 2 and we employ $\sigma = 0.15$ for $64^2$-pixel images in our experiments. Moreover $\sigma$ can be determined using only the real images.

*4.2. Comparison to conventional noise-space*

Given a constant $\sigma$, we can obtain the conventional noise-space as

$$\Phi_t^{\mathrm{n}}(y) := \Phi_{t-1}^{\mathrm{n}}(y) + \epsilon_t, \tag{4}$$

8

with $\epsilon_t \sim N(0, \sigma), t \in [T]$, and $\Phi_0^n(y) = y$. Equation (4) provides a discrete representation of data with additive noises. Figure 1 (a) shows an example of noise-space using a real image. We refer to GAN trained using the noise-space as Noise-Space (NS) GAN. Considering the sum of normal distributions, the noise-space defined by Eq.(4) is equivalent to

$$
\begin{aligned}
\Phi_t^n(y) &= y + \epsilon_1 + \epsilon_2 + \ldots + \epsilon_t, \\
&= y + \hat{\epsilon}(t),
\end{aligned}
\tag{5}
$$

where $\hat{\epsilon}(t) \sim N(0, \sigma \cdot t)$, $t \in [T]$. Therefore, the drawback of the noise-space is that it increases the variance of data, and makes the original data harder to be learned.

Note that the proposed noisy scale-space using Eq.(3) can be rewritten in a closed-form as

$$
\Phi_t(y) = \underbrace{k *^{(t)} y}_{\text{smoothed data}} + \underbrace{k *^{(t-1)} \epsilon_1 + k *^{(t-2)} \epsilon_2 + \ldots + k *^{(1)} \epsilon_{t-1} + k *^{(0)} \epsilon_t}_{\text{low-to high-frequency noises}},
\tag{6}
$$

where $k *^{(t)}$ denotes the $t$-times convolution with kernel $k$, and $k *^{(0)} \epsilon_t = \epsilon_t$. Thus, our noisy scale-space balances the variances of the data-term and the noise-terms using the smoothing. A finding is that we can keep the data variance using a constant $\sigma$ as demonstrated in Figure 2. Equation (6) also shows that the noisy scale-space replaces the high-frequency information in data by a set of noises with the corresponding frequencies. We demonstrate this property using a synthetic dataset in the following section.

Table 1: Fréchet inception distance (FID) and inception score (IS) for Celeba dataset by DCGAN trained using (1st, ..., 7th columns) two-step scale-space $\{t, 0\}$ with equal periods using $t = 0, 2, ..., 128$, respectively, and (8th column) a scale-space using exponential annealing with $t = 128$: $t = 0$ is the baseline. The mean of FID and IS were computed within 20 trials.

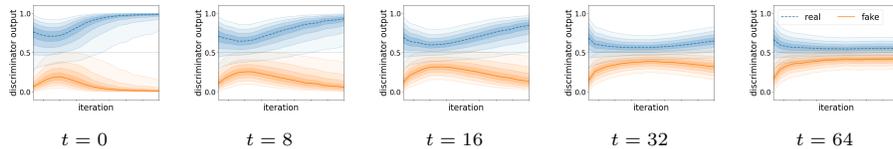| Smoothing | $t = 0$ | $t = 4$ | $t = 8$ | $t = 16$ | $t = 32$ | $t = 64$ | $t = 128$ | Annealing |
|---|---|---|---|---|---|---|---|---|
| FID ($\downarrow$) | **21.36** | 21.94 | 24.70 | 25.06 | 25.29 | 26.87 | 60.34 | 21.49 |
| IS ($\uparrow$) | 2.39 | 2.39 | 2.35 | 2.33 | 2.30 | 2.31 | 2.14 | 2.37 |

Figure 3: The prediction curves by discriminator for (blue doted line) the projected real data and (orange line) the projected fake data by DCGAN trained using CelebA dataset with (1st, ..., 5th columns) the fixed smoothing with $t$: $t = 0$ is the baseline. The percentiles of discriminator output were visualized in epoch wise.

### 4.3. Comparison to conventional scale-space

The conventional scale-space is a multi-scale representation of image in which the smoothing is applied to data recursively as

$$\Phi_t^{\mathrm{s}}(y) := k * \Phi_{t-1}^{\mathrm{s}}(y), \tag{7}$$

with $\Phi_0^{\mathrm{s}}(y) = y$ and $t \in [T]$. As shown in Figure 1 (b), the scale-space continuously removes the low-level information, e.g., textures and details, and provides the high-level information that is invariant to the scale.

We first re-examine the effect of diffusion in GAN optimization. Figure 3 visualizes the prediction curves of DCGAN [16] based on CelebA [39] where we trained the networks using the fixed smoothing time ($t$) applied to both real and fake data. We have chosen DCGAN since the convolution networks is the essential architecture of the modern GANs, e.g., [42, 43, 44, 45, 46, 47, 48]. As demonstrated in Figure 3, a larger $t$ results in better prediction curves where both $D(\Phi_t^{\mathrm{s}}(x))$ and $D(\Phi_t^{\mathrm{s}}(G(z)))$ converge to 0.5. Thus, the diffusion will make the problem more easy to learn by GANs. Note that the quality of $G(z)$ degenerates with $t$ since there is information loss due to the diffusion. Therefore, it is natural to anneal $t$ over the training process.

However, the conventional scale-space cannot improve the accuracy of GANs in general. Table 1 summaries the Fréchet inception distance (FID) [22] and the inception score (IS) [18] by DCGAN for CelebA [39] using scale-spaces, where we first used a fixed smoothing $t$ for 5 epochs then used the original data ($t = 0$) for 5 epochs. More details about the experimental set-up are summarized in

10

Section 5.1. As shown in Table 1, GANs using scale-spaces are often inferior to the baseline without smoothing.

Our key observation is that the scale-space has a side effect in the GAN optimization where the smoothed data (real and fake data) make the generator create a smoothed data and thus shrink the capability of the generator to learn high-frequency information in further training steps. To demonstrate this side effect, we propose to visualize the true probability distributions of both reals and fakes using a synthetic dataset that is created based on eight of vertical Hadamard bases ($\{B_i\}^8$) shown in Figure 4 (top). Given coefficients $\{\alpha_i\}^8$ for the eight basis, $\sum_i \alpha_i B_i$ produces an $8 \times 8$ pixel image as real data that looks a white/grey vertical stripes. The coefficients are the true probability distribution of real data. Given fake data, we can fit the bases to the fake image and compute the coefficients of eight bases that reflect the probability distribution of fakes with the fitting residuals.

We trained the basic GAN [1] based on the synthetic data using the smoothing with fixed $t = 8$. Figure 4 illustrates (top) the Hadamard bases, (middle) example of the original reals, smoothed-reals, smoothed-fakes, and the original fake data, and (bottom) the distribution of their Hadamard coefficients. There are two observations: (bottom-left) The coefficients of smoothed fakes $\Phi(G(z))$ follow the those of $\Phi(x)$ of which high-frequency coefficients are suppressed by the diffusion, i.e., the model learned the probability distribution of the smoothed data. (bottom-right) However, the diffusion decreased the diversity of high-frequency coefficients of fakes $(G(z))$. This means that the data smoothing reduces the overlap between the fake distribution with the distribution of reals with fine details that will be given in further steps.

We visualize the effect of the noisy scale-space using the synthetic dataset in Figure 5 in which we used a fixed $t = 8$. Figure 5 shows that (left) the projected data are smoothed yet (right) the diversity of high-frequency coefficients of fakes are preserved as expected.

Here, we have studied the side-effect of the diffusion. Interestingly, this kind of phenomenon due to a data transformation in GAN optimization is called
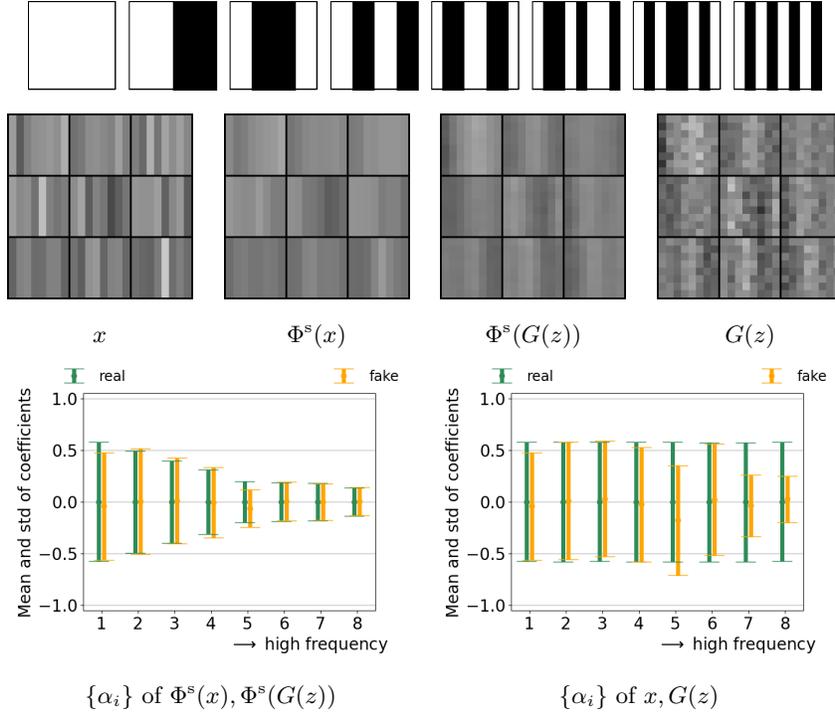
Figure 4: GAN with synthetic data using the scale-space: (top) The Hadamard bases $(B_1, ..., B_8)$, (middle part) nine examples of the original reals $(x)$, the smoothed reals $(\Phi^s(x))$ with fixed $t = 8$, the smoothed fakes $(\Phi^s(G(z)))$ with fixed $t = 8$, and the original fakes $(G(z))$, (bottom left) (y-axis) the mean and std. of the coefficients for (x-axis) the Hadamard bases $B_1, ..., B_8$ with low to high frequencies within the smoothed reals (green) and the smoothed fakes (orange), and (bottom right) those within the original reals and fakes. 200K of the synthetic images were created using $\alpha_i \sim U(-1, 1), \forall i$ with the uniform distribution $(U)$. We applied the smoothing to all data in the batch. The basic GAN [1] was trained using Adam with the learning-rate scale of $\eta = 2 \times 10^{-5}$ for 200 epochs.

'leaking' that has been studied on, e.g., rotation and hue change [17]. The contributions of our work are that we visualize the leaking of the smoothing using the synthetic data; and we propose the prescription to mitigate the side-effect of diffusion in GAN training.
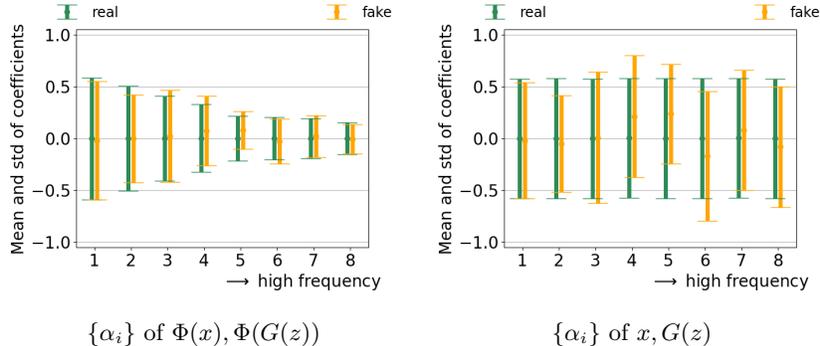
Figure 5: GAN with synthetic data using the proposed noisy scale-space: (left) The mean and std. of the Hadamard coefficients (y-axis) with low to high frequencies (x-axis) within the projected reals (green) and the projected fakes (orange), and (right) those within the original reals and fakes using the noisy scale-space with fixed $t = 8$.
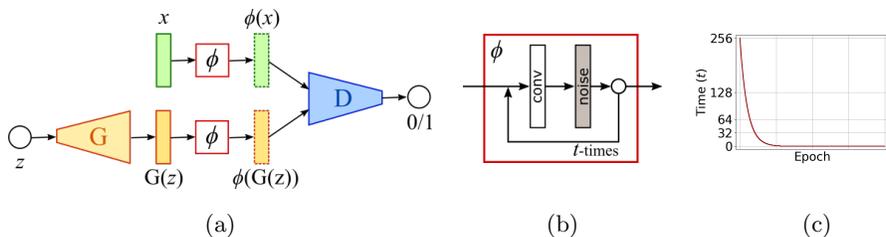


Figure 6: (a) The pipeline of GAN training with function $\Phi$ that is applied to the half of both real data $(x)$ and fake data $(G(z))$; (b) The module of the noise scale-space that consists of $t$-times repetition of the convolution-layer using the Gaussian kernel with the random noise-layer; and (c) The exponential annealing of $t$ using the power of $\beta = 20$ with $T = 256$.

### 4.4. Implementation of noisy scale-space

The data-based stabilization methods can be embedded into GANs as depicted in Figure 6 (a). We implement our data representation into DCGAN [16] that is the foundation of recent extensive studies, and call it Noisy Scale-Space (NSS) GAN. The function $\Phi$ consists of the repetition of the smoothing and the noise-injection layers (Figure 6b) that can be computed efficiently in parallel. Regarding the annealing of the time parameter $(t)$ that determines the magnitude of stabilization effect, we use an exponential function as

$$t^{(i)} := T \cdot \exp\left(i/\beta\right), \tag{8}$$

13

where $i \in [0, 1]$ is the relative iteration in the optimization process that starts at $i = 0$, $\beta$ is the power of decay, and $T$ is the initial time $t^{(0)}$. We use $\beta = 20$ in order to apply the data transformation in the early stage of GAN optimization. Figure 6 (c) shows the annealing curve of $t$ with $T = 256$. We have observed that the exponential function achieves better accuracy than others including the step function and the decaying-wave function. Also we employ an implementation technique [14] in which we apply the filtering function $\Phi$ to only half of data in real-and fake-batches for obtaining a stable and accurate generator.

## 5. Experimental results

In order to empirically demonstrate the expected stabilization effect of the noisy scale-space for GAN optimization, we conduct three experiments: an ablation study on the initial scale in data spaces, a primal experiment in which we compare the proposed NSS with the state-of-the-art GANs based on DCGAN, and an additional experiment using StyleGAN2.

### 5.1. Experimental set-up for primal experiments

In the preliminary experiment, we compare our NSS-GAN with potential competitors: DCGAN (GAN) as the baseline, GAN using the scale-space (SS-GAN), DCGAN using the noise-space (NS-GAN) where SS-GAN, NS-GAN, and NSS-GAN share the same architecture and the annealing function of $t$ except for the filtering function $\Phi$. Moreover, we employ GAN with the gradient regularization [3] (GAN-gr), LSGAN [49], WGAN-gp [18], and Dragan [20] as the state-of-the-arts in the second experiment.

We use four of the major datasets: CelebA [39], LSUN-Church, LSUN-Conference [40], and Oxford-Flowers [41] as the image generation tasks of faces, outdoor scenes, indoor scenes, and plants, respectively. CelebA consists of about 200K of celebrity face images. LSUN-Church and LSUN-Conference have about 126K images of outdoor scene of churches, and about 224K images of indoor scene of conferences, respectively. Oxford-Flowers has about 8K images of flowers. The images are resized to $64 \times 64$ pixels to meet the architecture of DCGAN.

14

We use the non-saturating loss that is known to be superior to the Wasserstein loss and others [3], and employ Adam [38] as one of the popular optimizers in GAN studies. For all the experiments, the mini-batch size is set to 128, and the number of training epochs is set to 10 epochs for CelebA, LSUN-Church and LSUN-Conference, and 100 epochs for Oxford-Flowers, respectively, based on their example sizes. The standard deviation of noise is set to $\sigma = 0.15$. The hyper-parameters that determine the quality of generated data are the initial time $T$ of SS, NS, and NSS-GANs, the learning-rate scale ($\eta$) and the first momentum coefficient ($b_1$) of Adam, and the regularization coefficient ($\lambda$) of GAN-gr. We conduct our experiments in two steps: a preliminary experiment on $T$ with fixed $\eta$ and $b_1$, and the final experiment using the tuned $\eta$, $b_1$, and $\lambda$ with fixed $T$. We perform each condition 20 individual times.

For quantitative evaluation of generated fake images, we use the Fréchet inception distance (FID) [22] with the inception score (IS) [18] that are widely used in GAN studies. FID measures the distance between the real data with fake data in feature space defined using a pre-trained network. IS measures the diversity of fakes in the feature space. Lower FID values with higher IS indicate better quality and diversity of fake data, respectively. We use FID as the primal metric that reflects the objective of GAN.

Table 2: The Fréchet inception distance (FID) over the initial time $T = 0, 32, 64, 128, 256$ for (column parts) CelebA, LSUN-Church, LSUN-Conference, and Oxford-Flowers by DCGAN using the scale-space (SS), the noise-space (NS), and our noisy scale-space (NSS): In order to demonstrate the stabilization effect over $T$, the learning-rate scale and the 1st momentum were fixed for each dataset such that the baseline DCGAN ($T = 0$) can be unstable.

|  | CelebA | | | Church | | | Conference | | | Flowers | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | SS | NS | NSS | SS | NS | NSS | SS | NS | NSS | SS | NS | NSS |
| $T=0$ | 24.60 | - | - | 70.94 | - | - | 127.32 | - | - | 201.02 | - | - |
| $T=32$ | 23.14 | 22.82 | 22.32 | 69.80 | 70.20 | 69.88 | 119.06 | 71.16 | 67.37 | 162.24 | 89.46 | 87.30 |
| $T=64$ | 23.16 | 23.01 | 22.43 | 69.05 | 69.60 | 69.67 | 98.73 | 72.21 | 70.33 | 120.57 | 93.46 | 88.68 |
| $T=128$ | 23.11 | 24.02 | 21.18 | 69.38 | 69.51 | 68.59 | 84.31 | 71.83 | 68.90 | 116.29 | 95.66 | 92.57 |
| $T=256$ | 22.31 | 24.29 | 21.79 | 70.58 | 69.61 | 67.38 | 87.46 | 73.11 | 70.47 | 100.84 | 95.47 | 94.24 |

Table 3: The tuned hyper-parameters of (columns) the experimented GANs for (rows) CelebA, LSUN-Church, LSUN-Conference, and Oxford-Flowers datasets: the learning-rate scale ($\eta \times 10^{-4}$), the 1st momentum ($b_1$) of Adam, and the regularization coefficient ($\lambda$) were selected using grid search based on the mean FID. GAN with the scale-space (SS), GAN with the noise-space (NS), and GAN with the noisy scale-space (NSS) share the fixed $T = 256$.

| | | GAN | GAN-gr | LSGAN | WGAN-gp | Dragan | SS | NS | NSS |
|---|---|---|---|---|---|---|---|---|---|
| CelebA | $\eta$ | 5 | 5 | 1 | 10 | 1 | 5 | 5 | 5 |
| | $b_1$ | 0.3 | 0.3 | 0.4 | 0.3 | 0.4 | 0.4 | 0.3 | 0.4 |
| | $\lambda$ | - | 1 | - | 20 | 10 | - | - | - |
| Church | $\eta$ | 2 | 5 | 0.5 | 10 | 1 | 5 | 10 | 5 |
| | $b_1$ | 0.3 | 0.3 | 0.5 | 0.3 | 0.4 | 0.3 | 0.3 | 0.3 |
| | $\lambda$ | - | 10 | - | 20 | 10 | - | - | - |
| Confer. | $\eta$ | 2 | 2 | 0.5 | 10 | 2 | 2 | 5 | 5 |
| | $b_1$ | 0.3 | 0.3 | 0.5 | 0.3 | 0.4 | 0.3 | 0.3 | 0.3 |
| | $\lambda$ | - | 5 | - | 20 | 10 | - | - | - |
| Flowers | $\eta$ | 2 | 5 | 0.5 | 10 | 2 | 5 | 5 | 5 |
| | $b_1$ | 0.4 | 0.4 | 0.5 | 0.3 | 0.5 | 0.3 | 0.3 | 0.3 |
| | $\lambda$ | - | 5 | - | 20 | 10 | - | - | - |

*5.2. Effect of initial time*

We first examine the initial time $T$ in Eq.(8) that determines the degree of the scale-space, the noise-space, and the proposed noisy scale-space. In order to observe the relative stabilization effect of the three data representations in comparison to the baseline GAN, we purposely use a condition that can make the baseline GAN unstable for each dataset. Concretely, we employ $(\eta, b_1) = (0.0002, 0.4)$ for CelebA and LSUN-Church, $(\eta, b_1) = (0.0005, 0.3)$ for LSUN-Conference, and $(\eta, b_1) = (0.0005, 0.4)$ for Oxford-Flowers datasets.

Table 2 shows the mean FID within the 20 trials over the initial time of $T = 0, 32, 64, 128, 256$ where $T = 0$ is the baseline DCGAN. Table 2 demonstrates that both the scale-space and the noise-space improved the quality of fakes compared to the baseline GAN; and the proposed NSS-GAN has achieved a better stabilization effect than SS-GAN and NS-GAN independent of the datasets and the initial time $T$. The scale ($t$) should start with a large value for any type of images and the initial value $T$ depends on the context and the scale of images essentially. In practice, we recommend $T = 256$ for $64^2$-pixel images and $T = 128$ for 512-pixel images and use them in the final experiments.

GAN         GAN-gr         LSGAN         WGAN-gp

Dragan         SS-GAN         NS-GAN         NSS-GAN

(a) CelebA

GAN         GAN-gr         LSGAN         WGAN-gp

Dragan         SS-GAN         NS-GAN         NSS-GAN

(b) Church

Figure 7: Fake images created by the generator of which FID is the closest to the mean of the individual trials for (top part) CelebA and (bottom part) LSUN-Church datasets using the baseline DCGAN (GAN), GAN with gradient regularization (GAN-gr) [3], LSGAN [49], WGAN-gp [18], (bottom) Dragan [20], and GANs using the scale-space (SS), the noise-space (NS), and the proposed noisy scale-space (NSS) with $T = 256$.

| GAN | GAN-gr | LSGAN | WGAN-gp |

| Dragan | SS-GAN | NS-GAN | NSS-GAN |

(a) Conference



| GAN | GAN-gr | LSGAN | WGAN-gp |

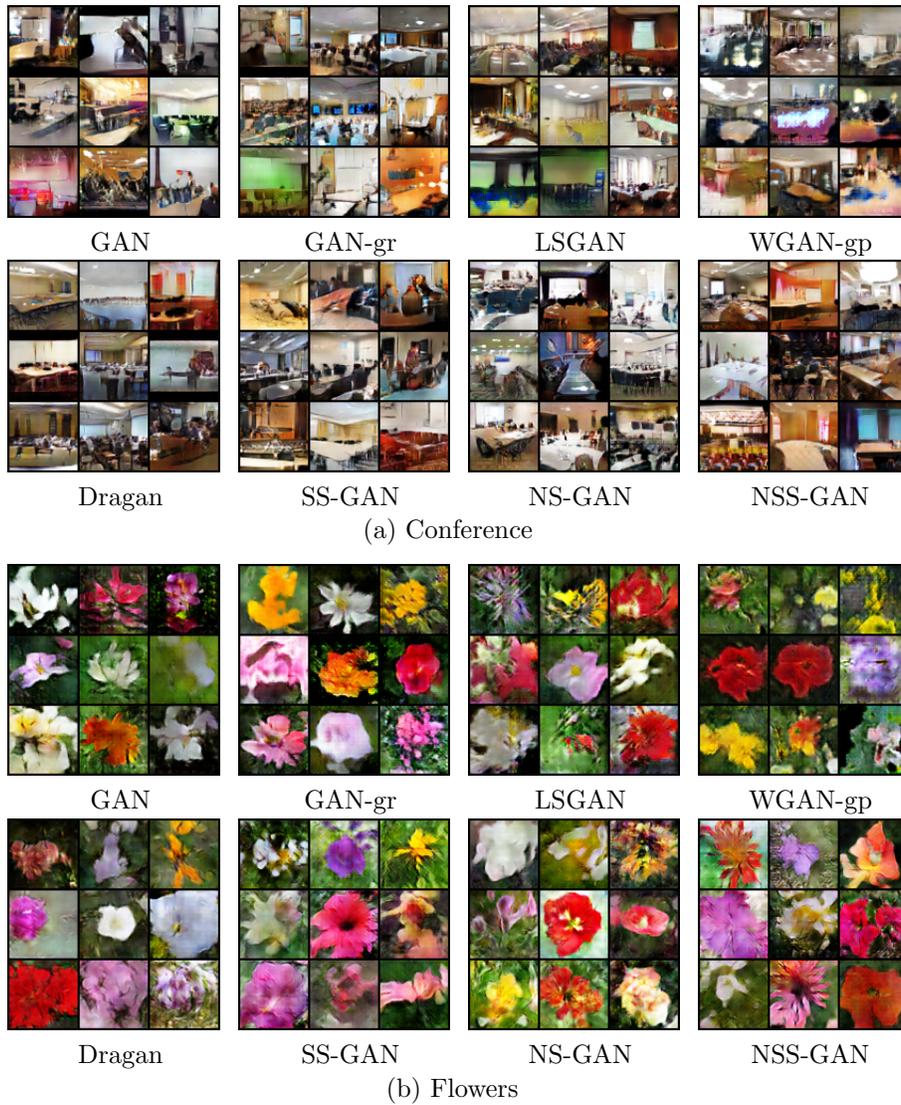| Dragan | SS-GAN | NS-GAN | NSS-GAN |

(b) Flowers

Figure 8: Fake images created by the generator of which FID is the closest to the mean of the individual trials for (top part) LSUN-Conference and (bottom part) Oxford-Flowers datasets using the baseline DCGAN (GAN), GAN with gradient regularization (GAN-gr) [3], LSGAN [49], WGAN-gp [18], (bottom) Dragan [20], and GANs using the scale-space (SS), the noise-space (NS), and the proposed noisy scale-space (NSS) with $T = 256$.

Table 4: The mean and std. of Fréchet inception distance (FID) and inception score (IS) within 20 individual trials using CelebA, LSUN-Church, LSUN-Conference, and Oxford-Flowers datasets by (rows in each part) the baseline DCGAN (GAN), GAN with gradient regularization (GAN-gr) [3], LSGAN [49], WGAN-gp [18], Dragan [20], and GANs using the conventional scale-space (SS), the noise-space (NS), and the proposed noisy scale-space (NSS) with $T = 256$: The learning-rate and 1st momentum were tuned for each condition based on the mean FID.

| | CelebA | | Church | |
| | FID ($\downarrow$) | IS ($\uparrow$) | FID | IS |
|---|---|---|---|---|
| GAN | $21.36 \pm 1.71$ | $2.39 \pm 0.06$ | $69.05 \pm 4.65$ | $2.98 \pm 0.07$ |
| GAN-gr | $19.87 \pm 1.19$ | $2.35 \pm 0.05$ | $61.67 \pm 5.07$ | $2.97 \pm 0.11$ |
| LSGAN | $33.38 \pm 5.10$ | $2.28 \pm 0.06$ | $75.88 \pm 9.56$ | $3.02 \pm 0.09$ |
| WGAN-gp | $42.16 \pm 2.77$ | $2.45 \pm 0.07$ | $102.61 \pm 18.48$ | $2.90 \pm 0.14$ |
| Dragan | $20.36 \pm 1.25$ | $2.35 \pm 0.05$ | $\mathbf{50.41} \pm 2.68$ | $2.91 \pm 0.04$ |
| SS-GAN | $21.49 \pm 1.49$ | $2.37 \pm 0.05$ | $61.20 \pm 3.53$ | $2.97 \pm 0.08$ |
| NS-GAN | $20.50 \pm 1.14$ | $2.37 \pm 0.04$ | $60.02 \pm 4.47$ | $2.98 \pm 0.08$ |
| NSS-GAN | $\mathbf{19.45} \pm 1.15$ | $2.42 \pm 0.05$ | $58.86 \pm 3.92$ | $2.97 \pm 0.08$ |
| | Conference | | Flowers | |
| | FID | IS | FID | IS |
| GAN | $77.96 \pm 5.49$ | $4.19 \pm 0.11$ | $97.71 \pm 4.91$ | $2.98 \pm 0.09$ |
| GAN-gr | $70.82 \pm 3.80$ | $4.08 \pm 0.11$ | $99.25 \pm 5.59$ | $3.15 \pm 0.09$ |
| LSGAN | $81.90 \pm 9.49$ | $4.01 \pm 0.11$ | $121.81 \pm 8.50$ | $2.72 \pm 0.15$ |
| WGAN-gp | $121.05 \pm 7.60$ | $3.51 \pm 0.13$ | $129.05 \pm 7.70$ | $2.89 \pm 0.08$ |
| Dragan | $67.07 \pm 5.45$ | $4.18 \pm 0.16$ | $98.37 \pm 6.63$ | $3.00 \pm 0.09$ |
| SS-GAN | $80.07 \pm 5.44$ | $4.02 \pm 0.10$ | $91.93 \pm 5.12$ | $3.11 \pm 0.09$ |
| NS-GAN | $73.11 \pm 4.07$ | $4.02 \pm 0.08$ | $91.44 \pm 4.47$ | $3.07 \pm 0.08$ |
| NSS-GAN | $\mathbf{64.61} \pm 3.20$ | $4.21 \pm 0.14$ | $\mathbf{86.28} \pm 6.57$ | $3.09 \pm 0.10$ |

### 5.3. Comparison to state-of-the-arts

We now compare the proposed NSS-GAN with the state-of-the-arts of GAN-optimization: the baseline DCGAN (GAN), GAN with the gradient regularization [3], LSGAN [49], WGAN with gradient penalty (WGAN-gp) [18], Dragan [20], GANs using the scale-space (SS-GAN), and the noise-space (NS-GAN) using the four of datasets. Note that these GANs share the same backbone architecture of DCGAN. In order to make our result independent to the hyper parameters of Adam, we employ grid search and tuned the learning-rate scale ($\eta$) in combination with the first momentum coefficient ($b_1$) for each pair of model and dataset, while the second momentum of Adam was fixed as $b_2 = 0.999$ based on our pretest. For each model, we then choose the best condition using the mean FID within the 20 trials. Table 3 summarizes the tuned parameters in which the baseline GAN prefers a stable condition compared to the others.

Figure 7 and Figure 8 present fake images by the tested GANs with their tuned hyper-parameters, where we use the generator of which FID is the closest to the mean of the independent trials. Table 4 summarizes the mean and std. of FID and IS of the experimented GANs for the four datasets within the 20 trials, where the proposed NSS-GAN with the constant parameter of $T = 256$ has achieved better and comparable results than the state-of-the-arts, demonstrating the effectiveness of the presented NSS-GAN.

More importantly, Table 4 shows that the conventional scale-space (SS) GAN and the noise-space (NS) GAN were not consistently better than the baseline GAN. In contrast, the proposed NSS-GAN has consistently outperformed the baseline GAN, SS-GAN, and NS-GAN in FID. This indicates that the stabilization effect of the presented NSS-GAN is not the simple summation of those by the scale-space with the noise-space but due to the better use of their mutually complementary relationship in the GAN optimization.

### 5.4. Comparison to StyleGAN2-Ada

As an additional study, we apply the noise scale-space to StyleGAN2-Ada [17] and compare it with the original algorithm using MetFaces dataset [17] and

Table 5: The hyper-parameters of (left) StyleGAN2-Ada and (right) StyleGAN2 with the proposed noisy scale-sapce: the R1 regularization ($\gamma$) of StyleGAN2 and the noise std of the proposed method ($\sigma$). We also employed the initial value $T = 128$ for our noisy scale-space.

|  |  | StyleGAN2-Ada | StyleGAN2-NSS |
|---|---|---|---|
| MetFaces-512 | $\gamma$ | 1.64 | 0.82 |
|  | $\sigma$ | - | 0.05 |
| FFHQ-512 | $\gamma$ | 1.64 | 1.64 |
|  | $\sigma$ | - | 0.1 |

Flickr-Faces-HQ (FFHQ) dataset [48]. StyleGAN2-Ada [17] is a state-of-the-art generative model for high-resolution images using a variety of data augmentations. We denote our implementation by StyleGAN2-NSS.

StyleGAN2 series require 25Kimg-iterations to converge with images of the size $1024 \times 1024$. However, this requires a large-scale GPUs with a huge computational time. For the sake of reproducibility, we use down-scaled images of the size $512 \times 512$ that we call MetFaces-512 and FFHQ-512, respectively, and we also limit the training iterations to 5Kimg that is known to achieve reasonable results [17].

For FFHQ-512 dataset, we implement our StyleGAN2-NSS by replacing the augmentation term (i.e., Ada) [17] by the proposed NSS function. However, we have observed that both the vanilla StyleGAN2 without Ada and StyleGAN2-NSS can explode when applied to MetFaces-512. This means that the backbone network of StyleGAN2 does not work with MetFaces-512 without the Ada term. Thus, our implementation for MetFaces-512 includes the Ada term in combination with our NSS function.

We follow the experimental set-up, the evaluation metrics, and the hyper-parameters of the official StyleGAN2 implementation, including the R1 regularization weight ($\gamma$) that we have tuned within $\gamma = 0.82, 1.64, 3.28$ as recommended in [17]. We tune the noise std of NSS based on the image variance curve of the images and also employ $T = 128$ as the initial condition. Table 5 summarizes the tuned hyper-parameters.

Table 6 summarizes FID and IS by StyleGAN2-Ada and StyleGAN2 using the proposed noisy scale-space. Figure 9 visualizes fake images created by the

generator with the average FID within the individual trials. Table 6 and Figure 9 show that the presented NSS-GAN has successfully improved the accuracy and quality of generated images using the state-of-art StyleGAN2.

It has been reported[17] that StyleGAN2-Ada outperforms Progressive-GAN [28] in accuracy. Therefore, our experimental results implicitly demonstrate that StyleGAN2-NSS is superior to Progressive-GAN. Moreover, our noisy scale-space is independent of the architecture in contrast to Progressive-GAN that requires the hierarchical architecture of networks.

Table 6: Fréchet inception distance (FID-50K) and inception score (IS-50K) for (top part) MetFaces-512 and (bottm part) FFHQ-512 datasets by (left) StyleGAN2-Ada and (right) StyleGAN2 with the proposed noisy scale-space. The models were trained 5KImg-iterations. The mean and std of the metrics were computed within the individual 5 trials.

|  |  | StyleGAN2-Ada | StyleGAN2-NSS (ours) |
|---|---|---|---|
| MetFaces-512 | FID ($\downarrow$) | $18.30 \pm 1.87$ | $\mathbf{17.25} \pm 0.56$ |
|  | IS ($\uparrow$) | $3.79 \pm 0.14$ | $\mathbf{3.87} \pm 0.05$ |
| FFHQ-512 | FID ($\downarrow$) | $7.27 \pm 0.25$ | $\mathbf{5.52} \pm 0.12$ |
|  | IS ($\uparrow$) | $4.61 \pm 0.03$ | $\mathbf{4.91} \pm 0.08$ |

## 6. Conclusion

In the consideration of data manipulation for the stable optimization in GANs, we have proposed a discrete representation of data, called noisy scale-space (NSS), that gradually removes high-frequency information in image while adding noise, leading to a coarse-to-fine training of GANs. In order to observe the side-effect of the conventional scale-space in GAN optimization, we have proposed the synthetic dataset based on the Hadamard bases that visualizes the true distribution of the real and fake data. We have experimented with the proposed NSS using two backbone networks: DCGAN and StyleGAN2 based on the major datasets for natural image generation tasks. The experimental results have successfully demonstrated that: NSS-based GANs overtook the potential competitors and the state-of-the-arts in most cases.

22

StyleGAN2-Ada                                   StyleGAN2-NSS (ours)

Figure 9: Fake images based on (top) MetFaces-512 and (bottom) FFHQ-512 datasets by (left) StyleGAN2-Ada and (right) StyleGAN2 with the proposed noisy scale-sapce trained 5KImg-iterations. The generators with the average FID within the 5 trials were used.

A limitation of our method is that we assume the diffusion can simplify (the real) data. Concretely, our NSS-GAN is inferior to the original GAN when using MNIST [50] images that consist of 0/1 binary values. Obviously, smoothing the binary data increases the diversity of pixel values. Our assumption holds for natural images and our method yields the sufficient stabilization effect for GAN optimization irrespective of the content of images.

## References

[1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, Advances in neural information processing systems 27 (2014) 2672–2680.

[2] V. Nagarajan, J. Z. Kolter, Gradient descent gan optimization is locally stable, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 30, Curran Associates, Inc., 2017.

[3] L. Mescheder, A. Geiger, S. Nowozin, Which training methods for gans do actually converge?, in: International conference on machine learning, PMLR, 2018, pp. 3481–3490.

[4] H. Berard, G. Gidel, A. Almahairi, P. Vincent, S. Lacoste-Julien, A closer look at the optimization landscapes of generative adversarial networks, in: 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020.

[5] K. Kurach, M. Lučić, X. Zhai, M. Michalski, S. Gelly, A large-scale study on regularization and normalization in gans, in: International Conference on Machine Learning, PMLR, 2019, pp. 3581–3590.

[6] D. Wang, X. Qin, F. Song, L. Cheng, Stabilizing training of generative adversarial nets via langevin stein variational gradient descent, IEEE Transactions on Neural Networks and Learning Systems (2020) 1–13`doi: 10.1109/TNNLS.2020.3045082`.

[7] M. Arjovsky, L. Bottou, Towards principled methods for training generative adversarial networks, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, 2017.

[8] A. Krogh, J. A. Hertz, A simple weight decay can improve generalization, in: Advances in neural information processing systems, 1992, pp. 950–957.

[9] R. Sutton, Two problems with back propagation and other steepest descent learning procedures for networks, in: Proceedings of the Eighth Annual Conference of the Cognitive Science Society, 1986, pp. 823–832.

[10] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, S. Paul Smolley, Least squares generative adversarial networks, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2794–2802.

[11] S. Nowozin, B. Cseke, R. Tomioka, f-gan: Training generative neural samplers using variational divergence minimization, in: D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 29, Curran Associates, Inc., 2016, pp. 271–279.

[12] L. Cai, Y. Chen, N. Cai, W. Cheng, H. Wang, Utilizing amari-alpha divergence to stabilize the training of generative adversarial networks, Entropy 22 (4) (2020) 410.

[13] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein generative adversarial networks, in: International conference on machine learning, PMLR, 2017, pp. 214–223.

[14] S. Jenni, P. Favaro, On stabilizing generative adversarial training with noise, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 12145–12153.

[15] W. J. Townsend, M. A. Thornton, Walsh spectrum computations using cayley graphs, in: Proceedings of the 44th IEEE 2001 Midwest Symposium on Circuits and Systems. MWSCAS 2001 (Cat. No. 01CH37257), Vol. 1, IEEE, 2001, pp. 110–113.

[16] A. Radford, L. Metz, S. Chintala, Unsupervised representation learning with deep convolutional generative adversarial networks, arXiv preprint arXiv:1511.06434.

[17] T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, T. Aila, Training generative adversarial networks with limited data, Advances in Neural Information Processing Systems 33 (2020) 12104–12114.

[18] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, X. Chen, X. Chen, Improved techniques for training gans, in: D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 29, Curran Associates, Inc., 2016, pp. 2234–2242.

[19] T. Miyato, T. Kataoka, M. Koyama, Y. Yoshida, Spectral normalization for generative adversarial networks, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, 2018.

[20] N. Kodali, J. Abernethy, J. Hays, Z. Kira, On convergence and stability of gans, arXiv preprint arXiv:1705.07215.

[21] K. Roth, A. Lucchi, S. Nowozin, T. Hofmann, Stabilizing training of generative adversarial networks through regularization, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 30, Curran Associates, Inc., 2017, pp. 2018–2028.

[22] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, S. Hochreiter, Gans trained by a two time-scale update rule converge to a local nash equilibrium, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 30, Curran Associates, Inc., 2017.

[23] F. Schaefer, H. Zheng, A. Anandkumar, Implicit competitive regularization in GANs, in: H. D. III, A. Singh (Eds.), Proceedings of the 37th International Conference on Machine Learning, Vol. 119 of Proceedings of Machine Learning Research, PMLR, 2020, pp. 8533–8544.

[24] D. Zhang, A. Khoreva, Progressive augmentation of gans, in: H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 32, Curran Associates, Inc., 2019, pp. 6249–6259.

[25] C. Zhang, S. Bengio, M. Hardt, B. Recht, O. Vinyals, Understanding deep learning requires rethinking generalization, arXiv preprint arXiv:1611.03530.

[26] T. Hazan, G. Papandreou, D. Tarlow, Adversarial perturbations of deep neural networks, Perturbations, Optimization, and Statistics, MITP (2017) 311–342.

[27] M. S. Sajjadi, G. Parascandolo, A. Mehrjou, B. Schölkopf, Tempered adversarial networks, in: International Conference on Machine Learning, PMLR, 2018, pp. 4451–4459.

[28] T. Karras, T. Aila, S. Laine, J. Lehtinen, Progressive growing of gans for improved quality, stability, and variation, arXiv preprint arXiv:1710.10196.

[29] A. Karnewar, O. Wang, Msg-gan: Multi-scale gradients for generative adversarial networks, in: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

[30] S. Zhao, Z. Liu, J. Lin, J.-Y. Zhu, S. Han, Differentiable augmentation for data-efficient gan training, Advances in Neural Information Processing Systems 33 (2020) 7559–7570.

[31] N.-T. Tran, V.-H. Tran, N.-B. Nguyen, T.-K. Nguyen, N.-M. Cheung, On data augmentation for gan training, IEEE Transactions on Image Processing 30 (2021) 1882–1897.

[32] H. Jun, R. Child, M. Chen, J. Schulman, A. Ramesh, A. Radford, I. Sutskever, Distribution augmentation for generative modeling, in: International Conference on Machine Learning, PMLR, 2020, pp. 5006–5019.

[33] T. Chen, X. Zhai, M. Ritter, M. Lucic, N. Houlsby, Self-supervised gans via auxiliary rotation loss, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 12154–12163.

[34] C. Sønderby, J. Caballero, L. Theis, W. Shi, F. Huszár, Amortised map inference for image super-resolution, in: International Conference on Learning Representations, 2017, pp. 1–17.

[35] A. P. Witkin, Scale-space filtering, in: Readings in Computer Vision, Elsevier, 1987, pp. 329–332.

[36] T. Lindeberg, Scale-space theory in computer vision, Vol. 256, Springer Science & Business Media, 2013.

[37] A. Bora, E. Price, A. G. Dimakis, Ambientgan: Generative models from lossy measurements, in: International conference on learning representations, 2018, pp. 1–22.

[38] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980.

[39] Z. Liu, P. Luo, X. Wang, X. Tang, Deep learning face attributes in the wild, in: Proceedings of International Conference on Computer Vision (ICCV), 2015.

[40] F. Yu, Y. Zhang, S. Song, A. Seff, J. Xiao, Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop, arXiv preprint arXiv:1506.03365.

[41] M.-E. Nilsback, A. Zisserman, Automated flower classification over a large number of classes, in: Indian Conference on Computer Vision, Graphics and Image Processing, 2008.

[42] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, P. Abbeel, Infogan: Interpretable representation learning by information maximizing generative adversarial nets, 2016, pp. 1–9.

[43] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, A. A. Efros, Context encoders: Feature learning by inpainting, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2536–2544.

[44] A. Odena, C. Olah, J. Shlens, Conditional image synthesis with auxiliary classifier gans, in: International conference on machine learning, PMLR, 2017, pp. 2642–2651.

[45] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, D. N. Metaxas, Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 5907–5915.

[46] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, in: Computer Vision (ICCV), 2017 IEEE International Conference on, 2017.

[47] A. Brock, J. Donahue, K. Simonyan, Large scale GAN training for high fidelity natural image synthesis, in: International Conference on Learning Representations, 2019.

[48] T. Karras, S. Laine, T. Aila, A style-based generator architecture for generative adversarial networks, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2019, pp. 4401–4410.

[49] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, S. Paul Smolley, Least squares generative adversarial networks, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2794–2802.

[50] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, Proceedings of the IEEE 86 (11) (1998) 2278–2324.