# InSPECT: Invariant Spectral Features Preservation of Diffusion Models

Baohua Yan, Qingyuan Liu, Jennifer Kava, Xuan Di

Columbia University

by2348@columbia.edu, ql2505@columbia.edu, jk4944@columbia.edu, sharon.di@columbia.edu

## Abstract

*Modern diffusion models (DMs) have achieved state-of-the-art image generation. However, the fundamental design choice of diffusing data all the way to white noise and then reconstructing it leads to an extremely difficult and computationally intractable prediction task. To overcome this limitation, we propose **InSPECT (Invariant Spectral Feature-Preserving Diffusion Model)**, a novel diffusion model that keeps invariant spectral features during both the forward and backward processes. At the end of the forward process, the Fourier coefficients smoothly converge to a specified random noise, enabling features preservation while maintaining diversity and randomness. By preserving invariant features, InSPECT demonstrates enhanced visual diversity, faster convergence rate, and a smoother diffusion process. Experiments on CIFAR-10, Celeb-A, and LSUN demonstrate that InSPECT achieves on average a 39.23% reduction in FID and 45.80% improvement in IS against DDPM for 10K iterations under specified parameter settings, which demonstrates the significant advantages of preserving invariant features: achieving superior generation quality and diversity, while enhancing computational efficiency and enabling faster convergence rate. To the best of our knowledge, this is the first attempt to analyze and preserve invariant spectral features in diffusion models.*

## 1. Introduction

Diffusion models (DMs) have emerged as one of the most powerful classes of generative models, achieving remarkable success across diverse data modalities including image-[7, 8, 24, 27], video- [1, 2, 9, 29] and audio-generation [13, 16, 18, 28]. Foundational frameworks such as Denoising Diffusion Probabilistic Models (DDPMs) [8] progressively add noise to data following a forward process, until it becomes a white noise. A backward process then recovers the data by removing added noise. Building upon these developments, researchers further explore advances in DDPMs [15, 17, 21, 24], improving generation quality.
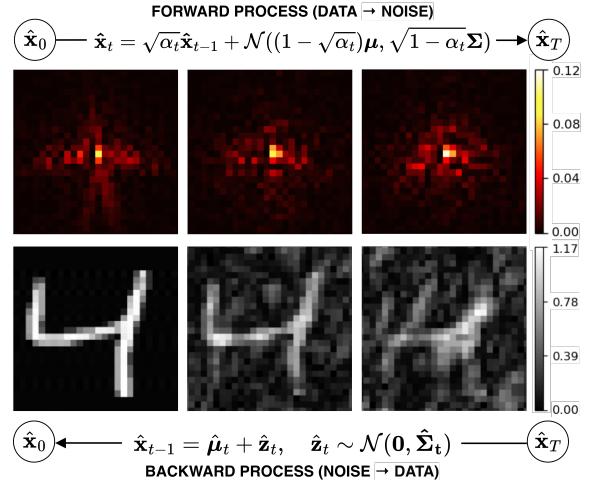


Figure 1. Visualization of InSPECT forward and backward processes. Input data $\hat{\mathbf{x}}_0$ is not completely destroyed; invariant component features are retained within the random noise at $\hat{\mathbf{x}}_T$

Despite these advances, current research works have yet to question a fundamental design choice of diffusion models: **must images be completely destroyed into white noise?** Notably, DDPMs-based models fully diffuse images into white noise during the forward process, erasing all useful information so that data reconstruction heavily depends on the accuracy of a learned denoising function, which is approximated by a U-NET neural network. While research works have been done on improving the U-NET architecture [21] or by implementing other architecture replacements [22], these approaches are generally time- and resource-intensive with a large degree of unpredictability, since the prediction of images or noises from nearly-white noises is extremely difficult and intractable. This inherent limitation motivates the development of a model that preserves informative features throughout the diffusion process, enabling a tractable, high-quality data generation while reducing computational costs.

Thus, a fundamental challenge lies in distinguishing invariant features from variant ones within image data. We

observe that Fourier transform holds great potential in this regard. By decomposing images (pixel space) via Fourier transform into Fourier coefficients (spectral space), we are able to calculate their statistical distributions, specifically the mean and variance. This is a good property that we empirically observe throughout image datasets, where images belonging to the same class exhibit a similar mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$ corresponding to their Fourier coefficients. Additionally, the Fourier coefficients contain invariant components with zero or near-zero variance, and we further note that these statistics differ across images of different classes (Figure 2). To examine if this spectral behavior is universal, we further analyze the standard deviation (STD) of Fourier coefficients across three widely used datasets (CIFAR-10 ($32 \times 32$) [14], Celeb-A ($64 \times 64$) [19], and LSUN ($128 \times 128$) [33]). Results in Table 1 shows a considerable number of spectral components exhibiting near-zero ($\leq 10^{-3}$) standard deviation, empirically confirming the existence of invariant spectral components.

Interestingly, the consistent occurrences of zero standard deviation is also observed in low-frequency components, generally related to overall geometric information, across natural images. We further visualize random noise $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ in spectral space by mapping its samples back to images. As illustrated in Figure 3, invariant spectral components contain geometric and class-related features, and therefore should be preserved during the diffusion process. These observations combined allow us to determine that the spectral space encodes class-dependent invariant spectral features. Notably, these statistical properties are not observable in pixel space. This establishes a consistent pattern of invariant spectral features across image datasets, which challenges the convection of complete diffusion to and reconstruction from white noise employed in standard diffusion models and inspires us to build diffusion models that intentionally preserve such invariant components during the generative process.

Motivated by the empirical observation that images have inherently invariant spectral features, we propose the **In**variant **Spect**ral Feature-Preserving Diffusion Model (**InSPECT**), a novel diffusion model based on DDPMs, but formulated and implemented in the spectral space. Beyond introducing the spectral space, we derive feature-preserving forward and backward processes that explicitly maintain the mean and variance of Fourier coefficients throughout the entire diffusion process, which allow the generative model to preserve invariant spectral features while adding random noise to diverse ones. By the end of the feature-preserving forward process, the Fourier coefficients of the initial images is destroyed into random noise $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, preserving invariant components while maintaining randomness and diversity. Since the noisiest Fourier components still contain invariant spectral features of the initial images, learning
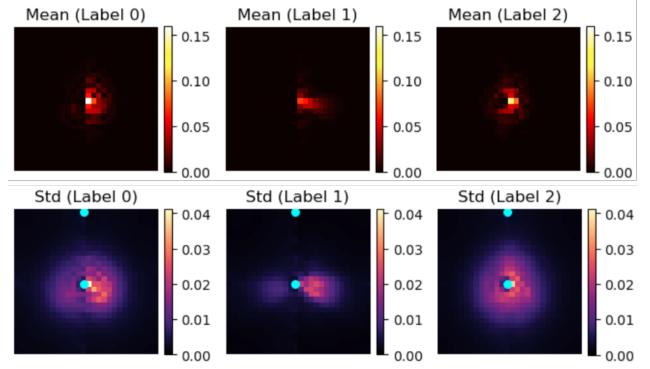


Figure 2. Heatmap illustration of Fourier coefficients mean $\boldsymbol{\mu}$ and standard deviation $\boldsymbol{\Sigma}^{1/2}$ in the spectral space of MNIST dataset corresponding to digits 0-2. Cyan points indicate invariant components, i.e. components with $\boldsymbol{\Sigma}^{1/2} = 0$.

Table 1. Empirical analysis of components with ($\leq 10^{-3}$) STD in CIFAR10, CelebA, and LSUN datasets. We discover that every image datasets contain a significant number of invariant or near-invariant components.

| | CIFAR10 | CelebA | LSUN |
|---|---|---|---|
| STD $= 0$ | 2 | 3 | 1 |
| STD $\leq 10^{-3}$ | $432 \pm 43$ | $9516 \pm 289$ | $14939 \pm 146$ |



Figure 3. Visualization of spectral random noise $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ corresponding to digits 0–2 in the MNIST dataset. We generate the random samples in spectral space and map them back to pixel space. Label-related features are preserved in the spectral space.

a denoising function is more tractable, achieving a balance between quality and efficiency in data generation.

By preserving invariant spectral features, InSPECT demonstrates enhanced visual diversity, faster convergence rate, and a smoother diffusion process. Experiments on three datasets (CIFAR-10 ($32 \times 32$), Celeb-A ($64 \times 64$), LSUN ($128 \times 128$)) show that InSPECT achieves $39.23\%$ FID improvement and $45.80\%$ IS improvement against the baseline diffusion model (DDPM) under specified parameter settings. These improvements demonstrate the significant advantages of preserving invariant features throughout the generative process: **achieving superior generation**

**quality and diversity with faster convergence rate**. Our key idea is summarized in Figure 1.

The main contributions of this paper are as follows:

(i) We discover and empirically confirm the existence of invariant spectral features within the spectral space;

(ii) we propose a feature-preserving model in spectral space, enabling the preservation of invariant spectral features to balance generation quality and efficiency;

(iii) we achieve both diverse and high-quality image generation with a faster convergence rate, presenting a computationally efficient alternative diffusion model to conventional DDPMs with promising initial results.

**Related works.** Many studies have examined diffusion processes in the frequency space [3, 6, 23]. These analysis consistently show that the forward noising schedule governs how different frequency components are activated and interact throughout the diffusion process [11, 31, 32]. Recent efforts have begun to address this gap from complementary angles. Jiralerspong et al. [10] investigate how shaping the inductive bias of the forward process in the frequency domain can better match the spectral statistics of natural data. EqualSNR [4] identifies that standard diffusion models corrupt different frequency bands at uneven rates, and proposes balancing the signal-to-noise ratio so that all frequencies receive more uniform treatment during diffusion.

## 2. Preliminaries

### 2.1. Denoising Diffusion Probabilistic Models

DDPMs consist of two main components: a forward (diffusing) process and a backward (denoising) process. The forward process [8] is defined by a Markov chain parameterized with Gaussian transitions that gradually add Gaussian noise $\epsilon_{\mathbf{t}}$ to progressively corrupt any image data $\mathbf{x}_0 \sim p_{data}(\mathbf{x})$, creating increasingly noisy data points. For any $t \in [0, T]$, the latent variable $\mathbf{x}_t$ is defined as:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_{\mathbf{t}} \qquad (1)$$

where $\epsilon_{\mathbf{t}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\bar{\alpha}_{\mathbf{t}} := \prod_{s=1}^{t} \alpha_s$. The denoising framework is then trained to invert this process by applying a reverse Markov chain starting from $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ with a noise predictor $\epsilon_\theta(\mathbf{x}_t, t)$:

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{1 - \beta_t}}\left[\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(\mathbf{x}_t, t)\right] + \sqrt{\beta_t}\mathbf{z} \quad (2)$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\beta_t := 1 - \alpha_t$. As noted in [8], using Bayes theorem, we calculate the posterior distribution of $\mathbf{x}_{t-1}$, denoted by $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ as defined below:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\tilde{\boldsymbol{\mu}}_t, \tilde{\beta}_t\mathbf{I}) \qquad (3)$$

where

$$\tilde{\boldsymbol{\mu}}_t := \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\mathbf{x}_0 \qquad (4)$$

$$\tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t \qquad (5)$$

### 2.2. Fourier Transform

Fourier transform projects a data point $\mathbf{x}_t(\mathbf{s}) \in C(\mathbb{R}^2)$ from the pixel space into the spectral (Fourier) space. The transformation decomposes the image into a weighted sum of trigonometric basis functions, producing spectral (Fourier) coefficients containing the amplitude and phase of each spatial frequency. Formally, the Fourier transform $\mathscr{F}$ is defined by an operator that maps input data $\mathbf{x}_t(\mathbf{s})$ to its corresponding $K$-dimensional Fourier coefficients, as shown below:

$$\begin{aligned} \mathscr{F} : C(\mathbb{R}^2) &\to \mathbb{C}^K \\ \mathbf{x}_t(\mathbf{s}) &\mapsto \hat{\mathbf{x}}_t \end{aligned} \qquad (6)$$

where $\mathscr{F}$ denotes the Fourier transform operator, and $\hat{\mathbf{x}}_t = (\hat{x}_{t,1}, \cdots, \hat{x}_{t,K})^T$ are multivariate Fourier coefficients. More specifically, the spectral composition, i.e. the Fourier transform $\mathscr{F}$ that connects the spatial variable $\mathbf{s} \in \mathbb{R}^2$ with the wavenumber $\mathbf{k}_j \in \mathbb{R}^2$ is shown through the following integral:

$$\hat{x}_{t,j} = \mathscr{F}[\mathbf{x}_t(\mathbf{s})] = \int_{\mathbb{R}^2} \mathbf{x}_t(\mathbf{s})\, e^{-i\mathbf{k}_j^T\mathbf{s}}\, d\mathbf{s}, \qquad (7)$$

where $e^{-i\mathbf{k}_j^T\mathbf{s}} = \cos(\mathbf{k}_j^T\mathbf{s}) - i\sin(\mathbf{k}_j^T\mathbf{s})$, $i$ is the imaginary unit and $\mathbf{k}_j$ represents the frequency of the trigonometric functions. Fourier coefficients contain structural information of low-frequency components (global, large-scale structure) and high-frequency components (local, small-scale details). This spectral representation provides a fully interpretable and mathematically tractable space in which statistical information can be calculated and preserved during the diffusion process.

The spectral decomposition of the input data $\mathbf{x}_t(\mathbf{s}) \in C(\mathbb{R}^2)$ is then defined by a linear combination of spatial basis functions and Fourier coefficients [25, 30]:

$$\mathbf{x}_t(\mathbf{s}) \approx \sum_{j=1}^{K} \hat{x}_{t,j} f_j(\mathbf{s}) \equiv \boldsymbol{f}^T(\mathbf{s})\hat{\mathbf{x}}_t \qquad (8)$$

where $f_j(\mathbf{s}) = \exp(i\mathbf{k}_j^T\mathbf{s})$ is the corresponding spatial basis function, with $\mathbf{s} \in \mathbb{R}^2$ representing the spatial variable and $\mathbf{k}_j$ representing the spatial wavenumbers (frequency vectors). Collectively, these form the Fourier basis, denoted by $\boldsymbol{f}(\mathbf{s}) = (f_1(\mathbf{s}), \cdots, f_K(\mathbf{s}))^T$, which spans the entire pixel space $C(\mathbb{R}^2)$. This bi-directional mapping using Fourier transform ensures that every dataset within the pixel space can be equivalently represented in an interpretable spectral space.

## 3. Methodology

In this section, we propose the **In**variant **Spect**ral Feature-Preserving Diffusion Model **(InSPECT)**, an extension of the Vanilla DDPM that introduces new forward and backward formulas defined in the spectral space. While the Vanilla DDPM completely destroys the image into white noise within the conventional pixel space, InSPECT indirectly destroys the image into random noise that preserves spectral features across the diffusion process by operating in the spectral space and obtaining a set of Fourier (spectral) coefficients corresponding to Fourier basis functions. Specifically, InSPECT utilizes Fourier transform which maps the diffusion process from pixel space into spectral (Fourier) space. For any initial data $\mathbf{x}_0$, we obtain its Fourier coefficient $\hat{\mathbf{x}}_0$ as follows:

$$\hat{\mathbf{x}}_0 = \mathscr{F}\mathbf{x}_0 \qquad (9)$$

where $\mathscr{F}$ is the Fourier transform in Eq (6). For each component of the Fourier coefficients $\hat{\mathbf{x}}_0$, we then calculate mean and variance across all sample points of the same class. These statistical properties are equivalently preserved by the new forward and backward processes we proposed, hence maintaining invariant spectral features during the diffusion process. Additionally, we introduce a frequency-weighted loss function that balances both low- and high-frequency predictions.

Ultimately, InSPECT is justified by its theoretical simplicity, relying only on statistical properties of Fourier coefficients, and its empirical effectiveness in preserving invariant features across the diffusion process, as demonstrated in subsequent experiments. The basic idea of InSPECT can also be transferred to other data modalities with an inherent spectral representation (e.g., audio, video, physical data or graphs), representing an interesting future direction.

### 3.1. Diffusion Process in InSPECT

Throughout this paper, the prior and posterior distributions are empirically determined to satisfy the general diffusion model assumption of a Gaussian distribution, which guarantees that $\hat{\mathbf{x}}_t$ at any time-step maintains Gaussian properties. We are hence able to analyze Fourier coefficients within the spectral space, and derive the forward and backward process formulas as follows:

### 3.1.1. Forward Process

We follow the forward noising paradigm of the Vanilla DDPM, but modify it to explicitly preserve the mean and variance by iteratively diffusing the data using noise with specified mean and variance at each time-step. Specifically, we compute the mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$ of $\hat{\mathbf{x}}_0$, the Fourier coefficients of the input images of the same class, and use these values to guide the forward process. This avoids complete destruction of invariant spectral features by directing the forward process towards a target random noise distribution, $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, that preserves these invariant features. Formally, we denote the mean of $\hat{\mathbf{x}}_0$ as $\boldsymbol{\mu} = (\mu_1, \cdots, \mu_K)^T$ and the variance of $\hat{\mathbf{x}}_0$ as $\boldsymbol{\Sigma} = diag\{\sigma_1, \cdots, \sigma_K\}$. We then define the single-step forward process $\hat{\mathbf{x}}_t$:

$$\hat{\mathbf{x}}_t = \sqrt{\alpha_t}\hat{\mathbf{x}}_{t-1} + \mathcal{N}((1-\sqrt{\alpha_t})\boldsymbol{\mu}, \sqrt{1-\alpha_t}\boldsymbol{\Sigma}) \qquad (10)$$

Given $\hat{\mathbf{x}}_0$, the closed-form (whole-step) forward process is defined by:

$$q(\hat{\mathbf{x}}_t|\hat{\mathbf{x}}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_0 + (1-\sqrt{\bar{\alpha}_t})\boldsymbol{\mu}, \ (1-\bar{\alpha}_t)\boldsymbol{\Sigma}) \quad (11)$$

We then find that at the final timestep $t = T$, when the noise-scheduler satisfies $\alpha_t = \bar{\alpha}_t = 0$, the distribution of the closed-form forward process, $q(\hat{\mathbf{x}}_T|\hat{\mathbf{x}}_0)$, simplifies to:

$$q(\hat{\mathbf{x}}_T|\hat{\mathbf{x}}_0) = \mathcal{N}(\boldsymbol{\mu}, \ \boldsymbol{\Sigma}) \qquad (12)$$

Thus, at the end of the forward process, the input data $\hat{\mathbf{x}}_0$ is diffused into a random noise that satisfies a Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}, \ \boldsymbol{\Sigma})$. This ensures that the resulting random noise maintains the same mean and variance as the original spectral representation of the images, successfully preserving invariant spectral features.

### 3.1.2. Backward Process

We continue building upon the backward (denoising) process of the Vanilla DDPM, but express it in the spectral space while leveraging the pixel-space denoising function. Let $m_\theta$ denote a neural network (e.g., a U-Net) operating in pixel space that is trained to predict the original image $\mathbf{x}_0$ from any noisy image $\mathbf{x}_t$, i.e.

$$\mathbf{x}_0 = m_\theta(\mathbf{x}_t, t) \qquad (13)$$

Since our diffusion model is defined in the spectral space, at each backward step for any Fourier coefficients $\mathbf{x}_t$, we first map back into pixel space by inverse Fourier transform $\mathbf{x}_t = \mathscr{F}^{-1}\hat{\mathbf{x}}_t$, apply the predictor $m_\theta$, and then map the predicted image $\mathbf{x}_0$ back into spectral space, as follows:

$$\hat{\mathbf{x}}_0 = \mathscr{F}m_\theta(\mathscr{F}^{-1}\hat{\mathbf{x}}_t, t) \qquad (14)$$

This design leverages the properties of both the spectral and pixel space: the tractable and interpretable spectral space for the preservation of invariant features, while the inductive bias of most denoising neural networks in pixel space allows for processing of spatially local features. Notably, the globally coupled property of spectral space makes denoising with $m_\theta$, which generally contains convolutional neural network, directly in spectral space inefficient and poorly conditioned. Hence, it is computationally and representationally advantageous to map $\hat{\mathbf{x}}_t$ back into pixel space
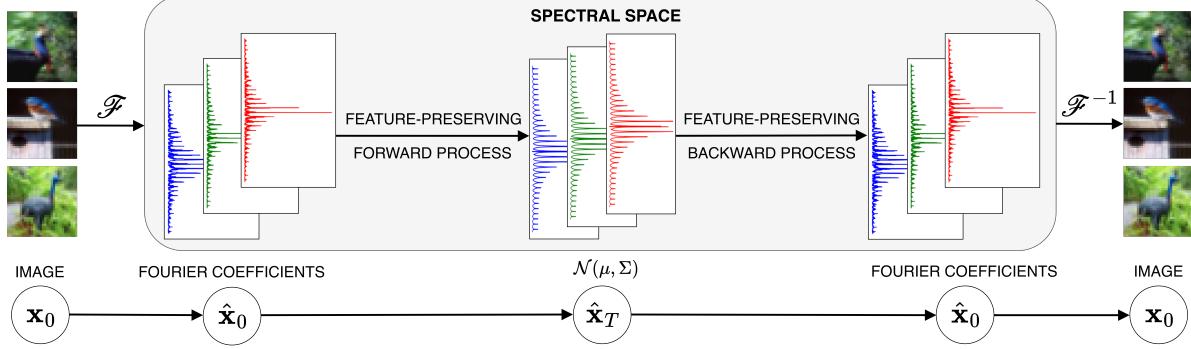
Figure 4. **The graphical model of Invariant Spectral Feature Preservation in Diffusion Models (InSPECT) considered in this work.** We convert a given image dataset, $\mathbf{x}_0$ into a Fourier coefficients and carry out the InSPECT forward and backward process, and convert the Fourier coefficients back into $\mathbf{x}_0$.

to process it through $m_\theta$, and then re-map the predicted image into spectral space.

Given the Fourier coefficients $\hat{\mathbf{x}}_t$ at any timestep $\mathbf{t}$, and the initial Fourier coefficients $\hat{\mathbf{x}}_0$, the posterior distribution $q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0, \hat{\mathbf{x}}_t)$ remains Gaussian. This follows from the fact that: (i) the forward process forms a Gaussian Markov chain, and (ii) the Fourier transform is an operator that preserves Gaussianity and do not alter variance structure, hence the posterior distribution necessarily remains Gaussian in spectral space. Moreover, we empirically observe that the Fourier coefficients of images are well-approximated by Gaussian distributions, so Gaussian assumption is consistent with the data and is not violated by operating within the spectral space. The posterior distribution of $\hat{\mathbf{x}}_{t-1}$ therefore takes the form of:

$$q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0, \hat{\mathbf{x}}_t) = \mathcal{N}(\hat{\boldsymbol{\mu}}_t, \hat{\beta}_t \boldsymbol{\Sigma}) \tag{15}$$

with mean

$$\hat{\boldsymbol{\mu}}_t = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \hat{\mathbf{x}}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \hat{\mathbf{x}}_0 + \gamma_t \boldsymbol{\mu} \tag{16}$$

The parameter $\gamma_t$ is defined to be:

$$\gamma_t := \frac{1 - \sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \beta_t - \frac{(1 - \bar{\alpha}_{t-1})(\sqrt{\alpha_t} - \alpha_t)}{1 - \bar{\alpha}_t} \tag{17}$$

Here, $\mu$ and $\Sigma$ represent the class-related mean and variance of the Fourier coefficients initially computed before the forward process. This enables a simplified model as the reverse Gaussian process can be expressed in closed form in terms of $\mu$, $\Sigma$ and the noise scheduler. Furthermore, since the mean $\hat{\boldsymbol{\mu}}_t$ is fully deterministic and uniquely determined by the model prediction of the current state, we can then simplify the one-step backward sampling as given by:

$$\hat{\mathbf{x}}_{t-1} = \hat{\boldsymbol{\mu}}_t + \hat{\mathbf{z}}_t, \quad \hat{\mathbf{z}}_t \sim \mathcal{N}(\mathbf{0}, \hat{\boldsymbol{\Sigma}}_{\mathbf{t}}) \tag{18}$$

This backward formula allows for a denoising process that is consistent with the forward process. The mean $\hat{\boldsymbol{\mu}}_t$ ensures that each sample converges back to the original distribution with a specific mean, while the Gaussian noise introduces stochasticity that maintains sample diversity and matches the variance of the original Fourier coefficients.

The joint use of pixel space and spectral space, together with the Gaussian property maintained throughout the diffusion process, greatly simplifies the backward process while preserving invariant spectral features. These properties allow InSPECT to achieve higher computational efficiency, faster convergence, and improved generation quality while enhancing sample diversity, as demonstrated in Section 4.

The forward and backward processes are summarized in Alg. 1 and Alg. 2.

### 3.2. Power-law of Fourier coefficients and weighted loss function

Empirically, the Fourier coefficients $\hat{\mathbf{x}}_t$ of natural images satisfy power scaling behavior and are heavy-tailed in many natural systems [5, 12, 20, 26]. The Fourier spectrum of Fourier transform is dominated by a power function and with a power-law decay with the form:

$$|\hat{\mathbf{x}}_{t,j}|^2 \propto \|\mathbf{k}_j\|^{-\gamma} \tag{19}$$

where $\mathbf{k}_j$ denotes the wavenumber associated with the frequency of the $j$-th Fourier basis function. A larger $\mathbf{k}_j$ corresponds to higher-frequency components, which capture the fine-grained details of images.

As illustrated in Figure 5, both the mean and standard deviation of the Fourier coefficients (e.g., in CIFAR-10 dataset) decay rapidly as $\mathbf{k}_j$ increases. One key observation from Figure 5 is that the Fourier coefficients corresponding to high-frequency terms (approximately $k > 50$) are nearly zero. This implies that an unweighted diffusion model natu-

rally focuses on reconstructing low-frequency components while neglecting high-frequency information. To mitigate this imbalance, we introduce a frequency-weighted loss function that rescales each Fourier component by using its standard deviation as weighting factors:

$$\mathcal{L}_t := \left\| \mathbf{\Sigma}^{-1/2} \left( \hat{\mathbf{x}}_0^* - \hat{\mathbf{x}}_0 \right) \right\|^2$$
$$= \frac{1}{K} \sum_{k=1}^{K} \frac{\left( \hat{\mathbf{x}}_{0,k}^* - \hat{\mathbf{x}}_{0,k} \right)^2}{\max\{\sigma_k, 10^{-3}\}} \quad (20)$$

The inverse standard deviation weighting amplifies the contribution of frequencies with small variance, which are primarily high-frequency components, thereby encouraging the model to preserve fine details, while the lower-bound acts as a threshold to prevent the weight from becoming arbitrarily large hence ensuring numerical stability. In practice, this encourages the model to allocate capacity to reconstruct fine details rather than focusing on only coarse structure, which leads to better image generation quality (e.g., sharper edges and richer textures).
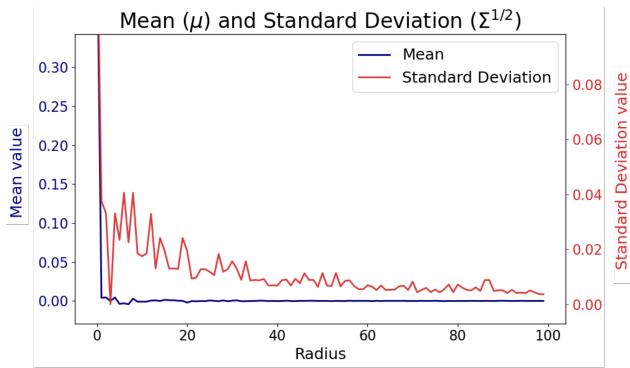


Figure 5. Mean and standard deviation of the CIFAR-10 dataset. $x$-axis is the frequency radius of the spectral coefficients. We pick red channel and show the curves related to one label.

## 3.3. Fourier transform in discrete pixel space

Although the Fourier transform is defined for functions in continuous space $C(\mathbb{R}^2)$, real-world practice (e.g., in scientific and engineering applications) typically involve real-valued data sampled on a discrete spatial grid, such as images. We assume that an image $\mathbf{x}_t$ in pixel space is defined on an $N_1 \times N_2$ regular spatial grid. The rectangular spatial mesh system is defined by a tensor product of two one-dimensional collocation sets:

$$\mathbb{S} = \mathbf{s}_1 \otimes \mathbf{s}_2 \quad (21)$$

where

$$\mathbf{s}_1 = (s_1^{(i)}; s_1^{(i)} = \frac{i}{N_1}, i = 0, \cdots, N_1 - 1)$$
$$\mathbf{s}_2 = (s_2^{(i)}; s_2^{(i)} = \frac{i}{N_2}, i = 0, \cdots, N_2 - 1) \quad (22)$$

Without loss of generality, $N_1$ and $N_2$ are assumed to be even numbers. Now we restrict our attention to two-dimensional discrete Fourier transform, which can be defined as

$$\hat{\mathbf{x}}_t(k_1, k_2) = \frac{1}{N_1 N_2} \sum_{n_1, n_2} \mathbf{x}_t(n_1, n_2) e^{-i(\frac{2\pi}{N_1} n_1 k_1 + \frac{2\pi}{N_2} n_2 k_2)}$$
$$= \frac{1}{N_1 N_2} \sum_{n_1, n_2} \mathbf{x}_t(n_1, n_2) e^{-i\omega_{n_1, k_1, n_2, k_2}}$$
$$= \frac{1}{N_1 N_2} \sum_{n_1, n_2} \mathbf{x}_t(n_1, n_2) \cos(\omega_{n_1, k_1, n_2, k_2})$$
$$- i \frac{1}{N_1 N_2} \sum_{n_1, n_2} \mathbf{x}_t(n_1, n_2) \sin(\omega_{n_1, k_1, n_2, k_2})$$
$$= \hat{\mathbf{x}}_t^{(R)}(k_1, k_2) - i\hat{\mathbf{x}}_t^{(I)}(k_1, k_2) \quad (23)$$

for $k_1 = -\frac{N_1}{2} + 1, -\frac{N_1}{2} + 2, \cdots, \frac{N_1}{2}$,
and $k_2 = -\frac{N_2}{2} + 1, -\frac{N_2}{2} + 2, \cdots, \frac{N_2}{2}$.

Since image $\mathbf{x}_t$ is real-valued, from rotational symmetry we have

$$|\hat{\mathbf{x}}_t(k_1, k_2)| = |\hat{\mathbf{x}}_t(-k_1, -k_2)|$$
$$\arg \hat{\mathbf{x}}_t(k_1, k_2) = \arg \hat{\mathbf{x}}_t(-k_1, -k_2) \quad (24)$$

where $\arg$ means the argument of complex numbers. Therefore, the dimension of $\hat{\mathbf{x}}_t(k_1, k_2)$ can be reduced to $N_1 \times \frac{N_2}{2}$ by removing half of the complex plane. Concatenating the real and imaginary parts, we can convert the Fourier coefficients into real-valued coefficients:

$$\hat{\mathbf{x}}_t(k_1, k_2) = \begin{cases} \hat{\mathbf{x}}_t^{(R)}(k_1, k_2), & k_2 \geq 0 \\ \hat{\mathbf{x}}_t^{(I)}(k_1, -k_2), & k_2 < 0 \end{cases} \quad (25)$$

which have the same dimension as $\mathbf{x}_t \in \mathbb{R}^{N_1 \times N_2}$.

## 4. Experiments

Given the theoretical formula of InSPECT, we look to empirically assessing how effective our proposed model is in practice. In particular, we aim to evaluate whether operating in spectral space indeed preserves invariant features by maintaining the mean and variance of Fourier coefficient throughout the forward and backward processes, and how well our model performs in the image generation process. Hence motivating the following research questions:

- **RQ1:** How does InSPECT perform image generation tasks compared to baseline diffusion models?

---

**Algorithm 1:** Training algorithm

---

**Input:** Data samples $S := \{\mathbf{x}_0^{(i)}\}_{i=1}^N$, number of
       diffusion steps $T$.

$\hat{\mathbf{x}}_0 \leftarrow \mathscr{F}\mathbf{x}_0$;

$\boldsymbol{\mu} \leftarrow \mathrm{Mean}(\hat{\mathbf{x}}_0)$;

$\boldsymbol{\Sigma} \leftarrow \mathrm{diag}\{\mathrm{Cov}(\hat{\mathbf{x}}_0)\}$;

**while** *not converged* **do**
    $t \sim \mathrm{Uniform}(\{1, 2, \cdots, T\})$;
    **Forward process:**
    $\hat{\mathbf{x}}_t = \sqrt{\alpha_t}\hat{\mathbf{x}}_{t-1} + \mathcal{N}((1-\sqrt{\alpha_t})\boldsymbol{\mu}, \sqrt{1-\alpha_t}\boldsymbol{\Sigma})$;
    **Predict sample:**
    $\hat{\mathbf{x}}_0^* = \mathscr{F}m_\theta(\mathscr{F}^{-1}\hat{\mathbf{x}}_t, t)$;
    **Compute loss:**
    $\mathcal{L}_t = \left\|\boldsymbol{\Sigma}^{-1/2}(\hat{\mathbf{x}}_0^* - \hat{\mathbf{x}}_0)\right\|^2$;
    Update $\theta$;

---

---

**Algorithm 2:** Sampling algorithm

---

**Input:** Random noise $\hat{\mathbf{x}}_T \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

**for** $t = T; t \geq 1; t = t - 1$ **do**
    $\hat{\mathbf{x}}_0 = \mathscr{F}m_\theta(\mathscr{F}^{-1}\hat{\mathbf{x}}_t, t)$ ;
    $\hat{\mathbf{x}}_{t-1} = \hat{\boldsymbol{\mu}}_t + \hat{\mathbf{z}}_t, \quad \hat{\mathbf{z}}_t \sim \mathcal{N}(\mathbf{0}, \tilde{\boldsymbol{\Sigma}}_{\mathbf{t}})$

**return** $\mathscr{F}^{-1}\hat{\mathbf{x}}_0$

---

- **RQ2:** Can InSPECT effectively preserve frequency invariance in the latent (spectral) space?
- **RQ3:** How does InSPECT compare to baseline diffusion models in terms of generation quality and diversity?

**Experiment Setting** Our experiments are conducted on three widely adopted datasets with different resolution: CIFAR-10 ($32 \times 32$) [14] and CelebA ($64 \times 64$) [19], and LSUN ($128 \times 128$) [33]. Furthermore, we adopted two standard metrics, Fréchet Inception Distance (FID) and Inception Score (IS), to quantitatively evaluate the generation quality and diversity. As a baseline, we compare our method, InSPECT, against the original diffusion model: Denoising Diffusion Probabilistic Models (DDPM) [8] with cosine noise schedule [21].

### 4.1. Performance on Unconditional Generation

We evaluate how the number of sampling steps and training iterations affects generation quality and diversity. A standard U-Net architecture [21] is employed for all diffusion models to ensure a fair comparison. We train all models under two training budgets: a resource-limited setting (2K iterations) and a longer training setting (10K iterations), and generate samples using 100 to 500 sampling steps. As shown in Table 2 in supplementary materials, **InSPECT consistently and often substantially outperforms**

**the baseline DDPM across nearly all metrics under both training budgets**, demonstrating its superior ability to preserve spectral structures while enhancing visual diversity. In particular, for the 10K iterations setting on the CIFAR-10 dataset, it achieves substantial gains; an average of 39.23% in FID and 45.80% in IS, indicating significantly better generation quality and diversity.

Specifically, Figure 7 further reports the FID and IS results for our InSPECT model and the baseline DDPM on the CIFAR-10 dataset. As illustrated in Figure 7, the FID of InSPECT rapidly decreases with more training iterations at $t = 500$, while the IS of InSPECT increases correspondingly. In contrast, both the FID and IS of DDPM show minimal improvement, even as the number of iterations increases from 2K to 10K. Notably, InSPECT (2K iterations) achieves higher IS and lower FID compared with the baseline model (10K iteration), showing its computational efficiency and fast convergent rate. These results indicate that InSPECT is able to preserve invariant spectral features, improve generation quality and diversity, while reducing computational cost and complexity.



(a) Inception Score (IS ↑) for different sampling steps under 2K and 10K training iterations.



(b) Fréchet Inception Distance (FID ↓) for different sampling steps under 2K and 10K training iterations.
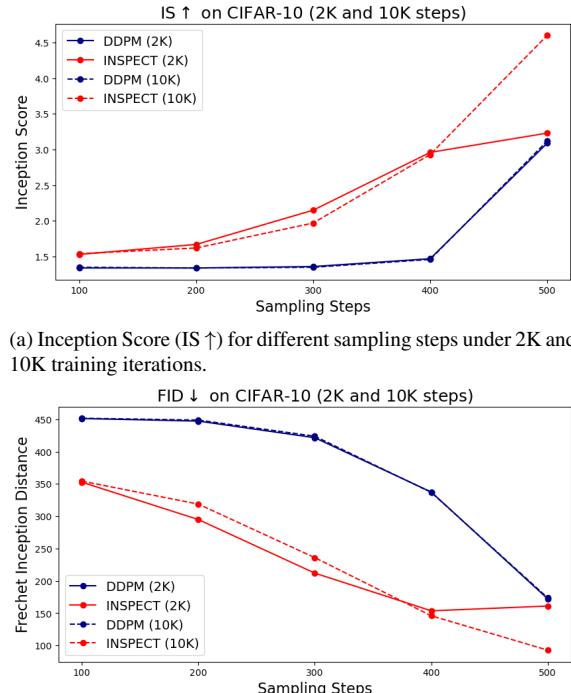
Figure 7. Comparison of DDPM and our proposed **InSPECT** model on CIFAR-10 across different sampling steps and training budgets. InSPECT achieves higher IS and lower FID than DDPM consistently across all sampling steps.
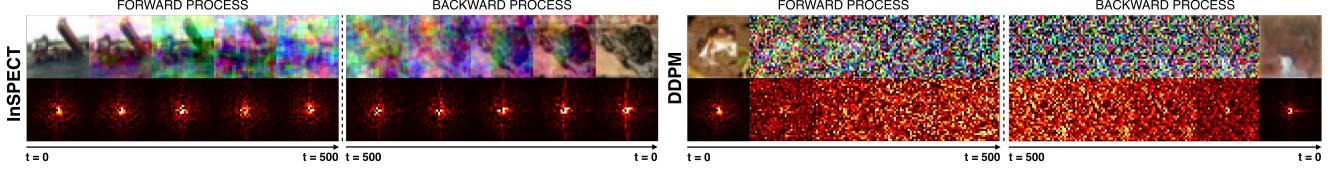
Figure 6. Visual comparison of forward and backward processes between InSPECT (**left**) and DDPM (**right**) by illustration of the images (**top row**) with their corresponding Fourier coefficients (**bottom row**) at different diffusion time-steps. InSPECT guides the image toward a specified random noise instead of a white noise, demonstrating significantly better generative quality.

## 4.2. Analysis of Frequency Invariance in Spectral Space

To further evaluate whether InSPECT can effectively preserve the spectral invariant features, we visualize the forward and backward processes of InSPECT and DDPM in pixel space and spectral space between $t = 0$ and $t = 500$. As illustrated in Figure 6, InSPECT preserves the invariant spectral features, as the distributions of the Fourier coefficients remain almost unchanged and evolve smoothly across diffusion steps, thereby avoiding the collapse of spectral structures while maintaining the diversity and randomness. In contrast, DDPM rapidly diffuses the image and fully destroys the spectral structures, leading to noisy Fourier coefficients and severe corruption of spectral features. By retaining the spectral features, InSPECT facilitates a more stable and smooth diffusion process, successfully reconstructs high-frequency coefficients (details of images), and ultimately achieves higher generation quality. These results highlight the key advantage of our approach: a feature-preserving diffusion process in spectral space, enabling improved generation quality and efficiency while maintaining generation diversity over standard diffusion models.
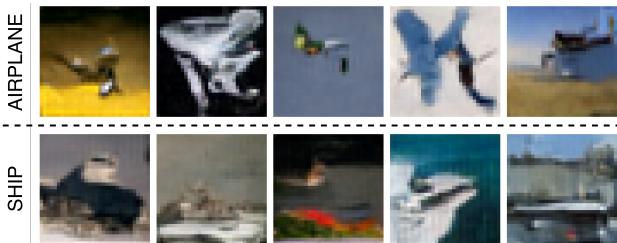
## 4.3. Performance on Conditional Generation



Figure 8. Class-conditional generation results on CIFAR-10. We specify two target classes (airplane and ship) and generate samples starting from class-related random noises. The generated images are well-aligned with the specified classes, demonstrating effective controllable class-conditional generation.

To evaluate the controllability of InSPECT, we perform class-conditional generation on the CIFAR-10 dataset with-

out retraining the entire model or training an external classifier model. During sampling, InSPECT starts from random noises with class-related mean and variance. Figure 8 shows representative samples generated for two specified classes (airplane and ship). InSPECT generates images that are consistent with the corresponding categories, demonstrating its effective class-conditional controllability. The results highlight that class-related spectral features can be preserved during our feature-preserving process, enabling more stable conditioning behavior and opening a promising direction for effective controllable generation based on invariant spectral features.

## 5. Conclusion

In this work, we discovered and analyzed the invariant features in the spectral space. Specifically, we empirically verified the existence of invariant spectral features. Building on this insight, we proposed InSPECT, a feature-preserving diffusion model that preserves the mean and variance of Fourier coefficients during the forward and backward processes. Extensive experiment on CIFAR-10, CelebA and LSUN datasets demonstrated that InSPECT not only maintains the spectral features and structures, but also effectively improves the generation quality and diversity. The supplementary materials provide discussions on additional experimental results and our limitations.

## References

[1] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 22563–22575. IEEE, 2023. 1

[2] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. 1

[3] Jonathan Crabbé, Nicolas Huynh, Jan Stanczuk, and Mihaela van der Schaar. Time series diffusion in the frequency do-

main. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024.* OpenReview.net, 2024. 3

[4] Fabian Falck, Teodora Pandeva, Kiarash Zahirnia, Rachel Lawrence, Richard E. Turner, Edward Meeds, Javier Zazo, and Sushrut Karmalkar. A fourier space perspective on diffusion models. *CoRR*, abs/2505.11278, 2025. 3

[5] David J Field. Relations between the statistics of natural images and the response properties of cortical cells. *Journal of the Optical Society of America A*, 4(12):2379–2394, 1987. 5

[6] Mathis Gerdes, Max Welling, and Miranda C. N. Cheng. GUD: generation with unified diffusion. *CoRR*, abs/2410.02667, 2024. 3

[7] Google DeepMind. Imagen 4. `https://deepmind.google/models/imagen/`, 2025. 1

[8] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 3, 7

[9] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *CoRR*, abs/2210.02303, 2022. 1

[10] Thomas Jiralerspong, Berton Earnshaw, Jason S. Hartford, Yoshua Bengio, and Luca Scimeca. Shaping inductive bias in diffusion models through frequency-based noise control. *CoRR*, abs/2502.10236, 2025. 3

[11] Diederik P. Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2021. Curran Associates Inc. 3

[12] Andrei Nikolaevich Kolmogorov. The local structure of turbulence in incompressible viscous fluid for very large reynolds numbers. *Proceedings of the Royal Society of London. Series A: Mathematical and Physical Sciences*, 434 (1890):9–13, 1991. 5

[13] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.* OpenReview.net, 2021. 1

[14] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, 2009. 2, 7

[15] Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent space. *arXiv*, 2022. DM Semantic Latent Space. 1

[16] Max W. Y. Lam, Jun Wang, Dan Su, and Dong Yu. BDDM: bilateral denoising diffusion models for fast and high-quality speech synthesis. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. 1

[17] Hang Li, Chengzhi Shen, Philip Torr, Volker Tresp, and Jindong Gu. Self-discovering interpretable diffusion latent directions for responsible text-to-image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024,* pages 12006–12016. IEEE, 2024. 1

[18] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo P. Mandic, Wenwu Wang, and Mark D. Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, pages 21450–21474. PMLR, 2023. 1

[19] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015. 2, 7

[20] Charlie Nash, Jacob Menick, Sander Dieleman, and Peter W. Battaglia. Generating images with sparse representations. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, pages 7958–7968. PMLR, 2021. 5

[21] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, pages 8162–8171. PMLR, 2021. 1, 7

[22] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 4172–4182. IEEE, 2023. 1

[23] Angus Phillips, Thomas Seror, Michael J. Hutchinson, Valentin De Bortoli, Arnaud Doucet, and Emile Mathieu. Spectral diffusion processes. *CoRR*, abs/2209.14125, 2022. 3

[24] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 10674–10685. IEEE, 2022. 1

[25] Fabio Sigrist, Hans R Künsch, and Werner A Stahel. Stochastic partial differential equation based modelling of large space–time data sets. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 77(1):3–33, 2015. 3

[26] Eero P Simoncelli and Bruno A Olshausen. Natural image statistics and neural representation. *Annual review of neuroscience*, 24(1):1193–1216, 2001. 5

[27] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 1

[28] Stability AI Introduces Stable Audio 2.5, the First Audio Model Built for Enterprise Sound Production at Scale. Stable audio 2.5. `https://stability.ai/news/stability-ai-introduces-stable-audio-25-the-first-audio-model-built-for-enterprise-sound-production-at-scale`, 2025. 1

[29] The Sora team. Sora 2. https://openai.com/index/sora-2/, 2025. 1

[30] Christopher K Wikle and Noel Cressie. A dimension-reduced approach to space-time kalman filtering. *Biometrika*, 86(4):815–829, 1999. 3

[31] Christopher Williams, Andrew Campbell, Arnaud Doucet, and Saifuddin Syed. Score-optimal diffusion schedules. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. 3

[32] Xingyi Yang, Daquan Zhou, Jiashi Feng, and Xinchao Wang. Diffusion probabilistic model made slim. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 22552–22562. IEEE, 2023. 3

[33] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 2, 7

# InSPECT: Invariant Spectral Features Preservation of Diffusion Models

## Supplementary Material

## 6. Theoretical proofs and discussions

### 6.1. Proof of backward posterior distribution

In this section, we will show the derivation of Eq 15. Recall that the forward process is:

$$\hat{\mathbf{x}}_t = \sqrt{\alpha_t}\hat{\mathbf{x}}_{t-1} + \mathcal{N}((1 - \sqrt{\alpha_t})\boldsymbol{\mu}, (1 - \alpha_t)\boldsymbol{\Sigma}) \tag{26}$$

Given $\hat{\mathbf{x}}_0$, the distribution of $\hat{\mathbf{x}}_t$ is:

$$q(\hat{\mathbf{x}}_t|\hat{\mathbf{x}}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\mu}, \ (1 - \bar{\alpha}_t)\boldsymbol{\Sigma}) \tag{27}$$

By Bayes' theorem, the posterior distribution can be written as:

$$q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0, \hat{\mathbf{x}}_t) \propto q(\hat{\mathbf{x}}_t|\hat{\mathbf{x}}_{t-1})q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0) \tag{28}$$

Note that $q(\hat{\mathbf{x}}_t|\hat{\mathbf{x}}_{t-1})$ and $q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0)$ are multidimensional Gaussian distributions, which can be expressed as:

$$
\begin{aligned}
q(\hat{\mathbf{x}}_t|\hat{\mathbf{x}}_{t-1}) &\propto \exp\left\{-\frac{1}{2(1-\alpha_t)\boldsymbol{\Sigma}}\|\hat{\mathbf{x}}_t - \sqrt{\alpha_t}\hat{\mathbf{x}}_{t-1} - (1 - \sqrt{\alpha_t})\boldsymbol{\mu}\|^2\right\} \\
&\propto \exp\left\{-\frac{1}{2(1-\alpha_t)\boldsymbol{\Sigma}}\left[\alpha_t\hat{\mathbf{x}}_{t-1}^T\hat{\mathbf{x}}_{t-1} - 2\sqrt{\alpha_t}\hat{\mathbf{x}}_{t-1}^T\hat{\mathbf{x}}_t + 2\sqrt{\alpha_t}(1 - \sqrt{\alpha_t})\hat{\mathbf{x}}_{t-1}^T\boldsymbol{\mu}\right]\right\}
\end{aligned} \tag{29}
$$

and

$$
\begin{aligned}
q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0) &\propto \exp\left\{-\frac{1}{2(1-\bar{\alpha}_{t-1})\boldsymbol{\Sigma}}\|\hat{\mathbf{x}}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\hat{\mathbf{x}}_0 - (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\mu}\|^2\right\} \\
&\propto \exp\left\{-\frac{1}{2(1-\bar{\alpha}_{t-1})\boldsymbol{\Sigma}}\left[\hat{\mathbf{x}}_{t-1}^T\hat{\mathbf{x}}_{t-1} - 2\sqrt{\bar{\alpha}_{t-1}}\hat{\mathbf{x}}_{t-1}^T\hat{\mathbf{x}}_0 - 2(1 - \sqrt{\bar{\alpha}_{t-1}})\hat{\mathbf{x}}_{t-1}^T\boldsymbol{\mu}\right]\right\}
\end{aligned} \tag{30}
$$

By substituting Eq 29 and Eq 30 into Eq 28, we have

$$
\begin{aligned}
q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0, \hat{\mathbf{x}}_t) \propto \exp\Bigg\{-\frac{1}{2\boldsymbol{\Sigma}}\Bigg[&\hat{\mathbf{x}}_{t-1}^T\left(\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}\right)I\hat{\mathbf{x}}_{t-1} \\
&- 2\hat{\mathbf{x}}_{t-1}^T\left((\sqrt{\alpha_t}\hat{\mathbf{x}}_t - \sqrt{\alpha_t}(1 - \sqrt{\alpha_t})\boldsymbol{\mu})\frac{1}{1-\alpha_t}\right. \\
&\left.+ (\sqrt{\bar{\alpha}_{t-1}}\hat{\mathbf{x}}_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\mu})\frac{1}{1-\bar{\alpha}_{t-1}}\right)\Bigg]\Bigg\}
\end{aligned} \tag{31}
$$

Let $\hat{\beta}_t^{-1} = \frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}$, then the posterior distribution can be simplified as:

$$q(\hat{\mathbf{x}}_{t-1}|\hat{\mathbf{x}}_0, \hat{\mathbf{x}}_t) = \mathcal{N}(\hat{\boldsymbol{\mu}}_t, \hat{\boldsymbol{\Sigma}}_t) \tag{32}$$

where

$$
\begin{aligned}
\hat{\boldsymbol{\Sigma}}_t &= \hat{\beta}_t\boldsymbol{\Sigma} \\
\hat{\boldsymbol{\mu}}_t &= \hat{\beta}_t\left((\sqrt{\alpha_t}\hat{\mathbf{x}}_t - \sqrt{\alpha_t}(1 - \sqrt{\alpha_t})\boldsymbol{\mu})\frac{1}{1-\alpha_t} + (\sqrt{\bar{\alpha}_{t-1}}\hat{\mathbf{x}}_0 + (1 - \sqrt{\bar{\alpha}_{t-1}})\boldsymbol{\mu})\frac{1}{1-\bar{\alpha}_{t-1}}\right) \\
&= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\hat{\mathbf{x}}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\hat{\mathbf{x}}_0 + \frac{1 - \sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t}\beta_t\boldsymbol{\mu} - \frac{(1 - \bar{\alpha}_{t-1})(\sqrt{\alpha_t} - \alpha_t)}{1 - \bar{\alpha}_t}\boldsymbol{\mu}
\end{aligned} \tag{33}
$$

## 6.2. Asymptotical discussion of mean value

In this section, we will discuss the asymptotical behavior of the mean value $\hat{\mathbf{x}}_t$ when $t$ is large. Consider a more general case where the forward process is defined as:

$$\hat{\mathbf{x}}_t = \sqrt{\alpha_t}\hat{\mathbf{x}}_{t-1} + \mathcal{N}(\lambda_t\boldsymbol{\mu}, (1-\alpha_t)\boldsymbol{\Sigma}) \tag{34}$$

where $\lambda_t$ is a time-dependent parameter. The closed-form forward process becomes:

$$q(\hat{\mathbf{x}}_t|\hat{\mathbf{x}}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_0 + \sqrt{\bar{\alpha}_t}\sum_{s=1}^{t}\frac{\lambda_t}{\sqrt{\bar{\alpha}_s}}\boldsymbol{\mu},\ (1-\bar{\alpha}_t)\boldsymbol{\Sigma}) \tag{35}$$

Note that if $\lambda_t = 1 - \sqrt{\alpha_t}$, then it reduces to our original formulation, i.e., Eq 11, whose mean value is:

$$\mathbb{E}[\hat{\mathbf{x}}_t] = \sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_0 + (1 - \sqrt{\bar{\alpha}_t})\boldsymbol{\mu} \tag{36}$$

When $t$ is large enough, $\bar{\alpha}_t \to 0$, thus $\mathbb{E}[\hat{\mathbf{x}}_t] \to \boldsymbol{\mu}$, which means that the mean value is convergent and bounded. Now we analyze the sufficient condition for the mean value to be convergent and bounded when $t$ is large. Consider $\sqrt{\bar{\alpha}_t}\sum_{s=1}^{t}\frac{\lambda_t}{\sqrt{\bar{\alpha}_s}}$:

$$
\begin{aligned}
\frac{\sqrt{\bar{\alpha}_{t+1}}\sum_{s=1}^{t+1}\frac{\lambda_{t+1}}{\sqrt{\bar{\alpha}_s}}}{\sqrt{\bar{\alpha}_t}\sum_{s=1}^{t}\frac{\lambda_t}{\sqrt{\bar{\alpha}_s}}} &= \frac{\sqrt{\bar{\alpha}_{t+1}}\left(\frac{\gamma_1}{\sqrt{\bar{\alpha}_1}} + \cdots + \frac{\gamma_{t+1}}{\sqrt{\bar{\alpha}_{t+1}}}\right)}{\sqrt{\bar{\alpha}_t}\left(\frac{\gamma_1}{\sqrt{\bar{\alpha}_1}} + \cdots + \frac{\gamma_t}{\sqrt{\bar{\alpha}_t}}\right)} \\
&= \sqrt{\alpha_{t+1}} + \frac{\sqrt{\alpha_{t+1}} \cdot \frac{\gamma_{t+1}}{\sqrt{\bar{\alpha}_{t+1}}}}{\frac{\gamma_1}{\sqrt{\bar{\alpha}_1}} + \cdots + \frac{\gamma_t}{\sqrt{\bar{\alpha}_t}}} \\
&\leq \sqrt{\alpha_{t+1}} + \frac{\sqrt{\alpha_{t+1}} \cdot \frac{\gamma_{t+1}}{\sqrt{\bar{\alpha}_{t+1}}}}{\frac{\gamma_t}{\sqrt{\bar{\alpha}_t}}} \\
&= \sqrt{\alpha_{t+1}} + \frac{\gamma_{t+1}}{\gamma_t}
\end{aligned} \tag{37}
$$

If there exists a constant $p < 1$ such that $\frac{\gamma_{t+1}}{\gamma_t} \leq p < 1$, then we have

$$\lim_{t\to\infty}\frac{\sqrt{\bar{\alpha}_{t+1}}\sum_{s=1}^{t+1}\frac{\lambda_{t+1}}{\sqrt{\bar{\alpha}_s}}}{\sqrt{\bar{\alpha}_t}\sum_{s=1}^{t}\frac{\lambda_t}{\sqrt{\bar{\alpha}_s}}} \leq q < 1 \tag{38}$$

By the ratio test, we conclude that $\sqrt{\bar{\alpha}_t}\sum_{s=1}^{t}\frac{\lambda_t}{\sqrt{\bar{\alpha}_s}}$ is convergent, and thus the mean value $\mathbb{E}[\hat{\mathbf{x}}_t]$ is convergent and bounded when $t$ is large. In summary, a sufficient condition for the mean value to be convergent and bounded when $t$ is large is that there exists a constant $p < 1$ such that $\frac{\lambda_{t+1}}{\lambda_t} \leq p < 1$.

## 7. Experimental Results

### 7.1. Performance on Unconditional Generation

In this section, we provide more experimental results on unconditional image generation tasks. We conduct experiments on three image datasets: CIFAR-10 ($32 \times 32$), CelebA ($64 \times 64$), and LSUN ($128 \times 128$). The quantitative comparison results between DDPM and our proposed InSPECT model are summarized in Table 2, Table 3, and Table 4, respectively. From these results, we can observe that our InSPECT model consistently outperforms the DDPM baseline across different datasets and sampling steps, demonstrating its effectiveness in enhancing image generation quality and diversity.

Table 2. Quantitative comparison of image generation quality between DDPM and our proposed InSPECT model on CIFAR-10 ($32 \times 32$) dataset. Results are reported as FID ($\downarrow$)/IS ($\uparrow$) for varying numbers of sampling steps $T$. Lower FID and higher IS denote better image quality and diversity.

| $T$ | 2K ITERATIONS | | 10K ITERATIONS | |
|---|---|---|---|---|
| | DDPM | InSPECT | DDPM | InSPECT |
| 100 | 451.20/1.34 | 352.68/1.53 | 451.65/1.35 | 354.37/1.54 |
| 200 | 447.26/1.34 | 294.93/1.67 | 448.84/1.34 | 318.78/1.62 |
| 300 | 421.53/1.36 | 211.98/2.15 | 423.84/1.35 | 235.97/1.97 |
| 400 | 337.50/1.47 | 153.62/2.96 | 336.90/1.46 | 145.96/2.93 |
| 500 | 172.21/3.09 | 161.00/3.23 | 173.79/3.12 | 92.76/4.60 |

Table 3. Quantitative comparison of image generation quality between DDPM and our proposed InSPECT model on CelebA ($64 \times 64$) dataset. Results are reported as FID ($\downarrow$)/IS ($\uparrow$) for varying numbers of sampling steps $T$.

| $T$ | 2K ITERATIONS | | 10K ITERATIONS | |
|---|---|---|---|---|
| | DDPM (2K) | InSPECT (2K) | DDPM (10K) | InSPECT (10K) |
| 100 | 465.58/1.18 | 417.52/1.71 | 457.96/1.18 | 364.76/1.97 |
| 200 | 454.80/1.19 | 373.32/1.75 | 442.03/1.18 | 294.58/2.21 |
| 300 | 422.21/1.27 | 307.56/1.99 | 403.46/1.22 | 216.91/2.29 |
| 400 | 364.44/1.53 | 237.20/2.36 | 327.50/1.66 | 144.41/2.20 |
| 500 | 333.81/2.40 | 193.55/2.40 | 202.33/2.25 | 122.05/2.44 |

Table 4. Quantitative comparison of image generation quality between DDPM and our proposed InSPECT model on LSUN ($128 \times 128$) dataset. Results are reported as FID ($\downarrow$)/IS ($\uparrow$) for varying numbers of sampling steps $T$.

| $T$ | 10K ITERATIONS | |
|---|---|---|
| | DDPM (10K) | InSPECT (10K) |
| 100 | 426.30/1.13 | 386.47/1.36 |
| 200 | 408.60/1.19 | 384.67/1.50 |
| 300 | 373.76/1.40 | 297.84/1.74 |
| 400 | 289.78/2.32 | 212.30/2.23 |
| 500 | 151.48/3.04 | 139.31/3.76 |

## 7.2. Performance on Conditional Generation

Class-conditional samples are shown in Figure 8, and an extended set of examples is provided in Figure 9. We perform class-conditional generation on the CIFAR-10 dataset without retraining the entire model or training an external classifier model, where the sampling process starts from random noises with corresponding class-related mean and variance. Examples highlight that class-related spectral features can be preserved during our feature-preserving process.

(a) class: airplane

(b) class: automobile

(c) class: bird

(d) class: cat

(e) class: deer

(f) class: dog

(g) class: frog

(h) class: horse

(i) class: ship

Figure 9. Class-conditional image generation results for each CIFAR-10 ($32 \times 32$) class demonstrate that InSPECT consistently produces images that are visually and statistically consistent with their respective specified categories.