

# Laplacian-guided Entropy Model in Neural Codec with Blur-dissipated Synthesis

## Supplementary Material (Appendix)

### 1. Denoising Diffusion Models

Denoising diffusion models are hierarchical latent variable models which generate sample through gradually removing noise from a randomly sampled white noise vector. The training procedure is comprised of two processes: diffusion or forward and denoising or backward. Diffusion process destroy the clean image and convert it to an approximately pure Gaussian noise during  $T$  time steps. The learnable denoising process then reconstructs the data distribution from white noise by reversing the diffusion process.

**Diffusion Process:** The diffusion process [1] can be described as a Markov chain, wherein each step of the forward path is defined by a Gaussian transition kernel:

$$q(\mathbf{z}_t|\mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \alpha_{t|t-1}\mathbf{z}_{t-1}, \sigma_{t|t-1}^2\mathbf{I}), \quad (1)$$

where  $\alpha_{t|t-1} \in R^+$  governs the extent to which the previous latent is retained, while  $\sigma_{t|t-1} \in R^+$  regulates the magnitude of the added noise. The dimension of the latent variables  $\mathbf{z}_1, \dots, \mathbf{z}_T$  is the same as that of the data  $\mathbf{x}$  or  $\mathbf{z}_0$ . An important property of the forward process is that any desired step  $\mathbf{z}_t$  can be directly sampled from  $\mathbf{x}$  using a closed-form solution, without needing to compute preceding steps:

$$q(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\mathbf{z}_t; \alpha_t\mathbf{x}, \sigma_t^2\mathbf{I}), \quad (2)$$

where  $\alpha_{t|t-1} = \alpha_t/\alpha_{t-1}$  and  $\sigma_{t|t-1}^2 = \sigma_t^2 - \alpha_{t|t-1}^2\sigma_{t-1}^2$ . The pre-specified hyperparameters  $\alpha_t$  typically exhibit a monotonically decreasing pattern from 1 to 0, while  $\sigma_t$  monotonically increases from 0 to 1. This pattern leads to a gradual corruption of the input image by Gaussian noise as  $t$  increases, resulting in  $q(\mathbf{z}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

**Denoising Process:** The true denoising distribution, which is tractable when conditioned on  $\mathbf{x}$  [1], can be written:

$$q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_{t \rightarrow t-1}(\mathbf{x}, \mathbf{z}_t), \sigma_{t \rightarrow t-1}^2\mathbf{I}), \quad (3)$$

where the distribution parameters can be computed as:

$$\begin{aligned} \sigma_{t \rightarrow t-1} &= \sigma_{t|t-1}\sigma_{t-1}/\sigma_t \\ \boldsymbol{\mu}_{t \rightarrow t-1} &= (\alpha_{t|t-1}\sigma_{t-1}^2/\sigma_t^2)\mathbf{z}_t + (\alpha_{t-1}\sigma_{t|t-1}^2/\sigma_t^2)\mathbf{x} \end{aligned} \quad (4)$$

To generate data, the true denoising process can be estimated by a learned denoising distribution  $p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t) := q(\mathbf{z}_{t-1}|\mathbf{z}_t, \hat{\mathbf{x}} = \phi_\theta(\mathbf{z}_t, t))$ , where  $\hat{\mathbf{x}}$  is predicted from diffused sample  $\mathbf{z}_t$  using a neural network  $\phi_\theta$ . Similar to Eq. 4,  $p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t)$  can be expressed by the approximation  $\hat{\mathbf{x}}$ :

$$\begin{aligned} p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t) &= \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{z}_t, t), \sigma_{t \rightarrow t-1}^2\mathbf{I}), \\ &= \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_{t \rightarrow t-1}(\hat{\mathbf{x}}, \mathbf{z}_t), \sigma_{t \rightarrow t-1}^2\mathbf{I}). \end{aligned} \quad (5)$$

**Training Objective:** The likelihood function  $\log p_\theta(\mathbf{x})$  is challenging to compute directly for training the model. So, during training, its evidence lower bound is maximized ( $\text{ELBO} \leq \log p_\theta(\mathbf{x})$ ), which can be expressed as:

$$\begin{aligned} \text{ELBO} &= \mathbf{E}_q[-\overbrace{D_{KL}(q(\mathbf{z}_T|\mathbf{x})||p(\mathbf{z}_T))}^{L_T} + \overbrace{\log p_\theta(\mathbf{x}|\mathbf{z}_1)}^{L_0}] \\ &+ \sum_{t=2}^T -\overbrace{D_{KL}(q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})||p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t))}^{L_{t-1}}. \end{aligned} \quad (6)$$

Within a well-defined noise scheduling, both  $L_0$  and  $L_T$  tend to approach approximately 0 and remain constant. Therefore, for training the diffusion model, it becomes adequate to optimize the  $L_{t-1}$  term, which is equivalent to comparing the learnable denoising process with the true denoising distribution. As both of these distributions are Gaussian, the expressions for the KL divergences have closed-form solutions and can be written as follows:

$$L_{t-1} \propto \mathbf{E}_q[||\boldsymbol{\mu}_{t \rightarrow t-1} - \boldsymbol{\mu}_\theta(\mathbf{z}_t, t)||^2] = \mathbf{E}_q[||\mathbf{x} - \hat{\mathbf{x}}||^2]. \quad (7)$$

In above formulation, the neural network directly predicts  $\hat{\mathbf{x}}$ . However, [1] discovered that optimization becomes simpler by predicting Gaussian noise instead. Hence, if we express  $\mathbf{z}_t = \alpha_t\mathbf{x} + \sigma_t\boldsymbol{\epsilon}$ , then the neural network  $\phi_\theta$  generates  $\hat{\boldsymbol{\epsilon}} = \phi_\theta(\mathbf{z}_t, t)$ , resulting in:

$$\hat{\mathbf{x}} = (1/\alpha_t)\mathbf{z}_t - (\sigma_t/\alpha_t)\hat{\boldsymbol{\epsilon}}. \quad (8)$$

As demonstrated in [2], using this specific parameterization, the final loss is obtained as follows:

$$\mathbf{E}_{t, \mathbf{x}, \boldsymbol{\epsilon}}[||\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}||^2] = \mathbf{E}_{t, \mathbf{x}, \boldsymbol{\epsilon}}[||\boldsymbol{\epsilon} - \phi_\theta(\alpha_t\mathbf{x} + \sigma_t\boldsymbol{\epsilon}, t)||^2]. \quad (9)$$

## 2. Additional Details on Blurring Diffusion Model

**Heat Dissipation as Gaussian Diffusion:** The heat dissipation process or blurring [4] can be expressed as a type of Gaussian diffusion. First, the marginal distribution of any time step noisy latent  $\mathbf{z}_t$  can be defined as follows:

$$q(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\mathbf{z}_t; \mathbf{A}_t\mathbf{x}, \sigma^2\mathbf{I}), \quad (10)$$

where  $\mathbf{A}_t = \mathbf{V}\mathbf{D}_t\mathbf{V}^T$  represents the dissipation or blurring operation.  $\mathbf{V}^T$  contains orthogonal Discrete Cosine Transform (DCT) basis, while the diagonal matrix  $\mathbf{D}_t = \exp(-\Lambda\tau_t)$  corresponds to the exponentiation of a weighting matrix for the frequencies  $\Lambda$ .  $\Lambda$  contains squared frequencies  $\lambda_{n,m} = -\pi^2(n^2/W^2 + m^2/H^2)$ , where  $W$  and  $H$  are the width and height of the image, and  $n \in \{0, \dots, W-1\}$  and  $m \in \{0, \dots, H-1\}$ . According to Eq. 10, any latent state  $\mathbf{z}_t$  is created by introducing a constant level of noise to a progressively blurred data point. When we transform the variables using the following transformations:  $\mathbf{f}_t = \mathbf{V}^T\mathbf{z}_t$  and  $\mathbf{f}_x = \mathbf{V}^T\mathbf{x}$ , the Gaussian diffusion process can be formulated in frequency space:

$$\begin{aligned} q(\mathbf{V}^T\mathbf{z}_t|\mathbf{V}^T\mathbf{x}) &= \mathcal{N}(\mathbf{V}^T\mathbf{z}_t; \mathbf{V}^T\mathbf{A}_t\mathbf{x}, \sigma^2\mathbf{V}^T\mathbf{I}\mathbf{V}) \Leftrightarrow \\ q(\mathbf{f}_t|\mathbf{f}_x) &= \mathcal{N}(\mathbf{f}_t; \mathbf{D}_t\mathbf{f}_x, \sigma^2\mathbf{I}). \end{aligned} \quad (11)$$

If we define a vector  $\lambda$  containing the diagonal elements of  $\Lambda$ , we can express  $\mathbf{d}_t$  as  $\exp(-\lambda\tau_t)$ , which corresponds to the diagonal elements of the matrix  $\mathbf{D}_t$ . With this reinterpretation, the diffusion process in frequency space can be written as follows:

$$q(\mathbf{f}_t|\mathbf{f}_x) = \mathcal{N}(\mathbf{f}_t; \mathbf{d}_t \odot \mathbf{f}_x, \sigma^2\mathbf{I}), \quad (12)$$

where  $\odot$  denotes elementwise vector multiplication. Eq. 12 shows that the marginal distribution of  $\mathbf{f}_t$  can be decomposed into individual scalar elements  $f_t^{(i)}$ . Likewise, the learnable inverse heat dissipation model  $p_\theta(\mathbf{f}_{t-1}|\mathbf{f}_t)$  can also be decomposed in a fully factorized manner. As a result, we have the option to describe the heat dissipation process and its inverse using scalar representations for each dimension  $i$ :

$$\begin{aligned} q(f_t^{(i)}|f_x^{(i)}) &= \mathcal{N}(f_t^{(i)}; d_t^{(i)}u_x^{(i)}, \sigma^2) \Leftrightarrow \\ f_t^{(i)} &= d_t^{(i)}f_x^{(i)} + \sigma\epsilon, \text{ with } \epsilon \sim \mathcal{N}(0, 1). \end{aligned} \quad (13)$$

Eq. 13 can be identified as a particular case of the standard Gaussian diffusion process that operates in frequency space, i.e.,  $f_t^{(i)} = \alpha_t f_x^{(i)} + \sigma_t \epsilon$ , where  $\alpha_t = d_t^{(i)}$  and  $\sigma_t = \sigma$ . What distinguishes this type of diffusion process from the standard one is the utilization of distinct noise schedules,

denoted as  $\alpha_t$  and  $\sigma_t$ , for each scalar element of the latent variable  $\mathbf{f}_t$ . In other words, the noise applied in this process exhibits non-isotropic characteristics. It's worth noting that while the marginal variance  $\sigma$  is shared across all scalar elements  $f_t^{(i)}$ , the specific noise schedules provide individual adjustments for each element.

In heat dissipation models, the Markov process  $q(\mathbf{f}_t|\mathbf{f}_{t-1})$  can be defined, corresponding to their chosen marginal distribution  $q(\mathbf{f}_t|\mathbf{f}_x)$ . By establishing an equivalence with Gaussian diffusion, this process can be effectively described using the following formulation:

$$\begin{aligned} q(\mathbf{f}_t|\mathbf{f}_{t-1}) &= \mathcal{N}(\mathbf{f}_t; \alpha_{t|t-1}\mathbf{f}_{t-1}, \sigma_{t|t-1}^2\mathbf{I}), \\ \text{where } \alpha_t &= \mathbf{d}_t, \sigma_t^{(i)} = \sigma \Rightarrow \alpha_{t|t-1} = \frac{\mathbf{d}_t}{\mathbf{d}_{t-1}}, \\ &\Rightarrow \sigma_{t|t-1}^2 = (1 - (\frac{\mathbf{d}_t}{\mathbf{d}_{t-1}})^2)\sigma^2. \end{aligned} \quad (14)$$

When  $\mathbf{d}_t$  is designed to have smaller values for higher frequencies,  $\sigma_{t|t-1}$  will introduce greater noise to the higher frequencies at each timestep. This results in the heat dissipation model erasing information from those frequencies more rapidly compared to the standard diffusion process.

**Inverse Heat Dissipation:** Similar to the standard diffusion model [1], the analytical expression for the true inverse heat dissipation process is obtained and can be written as follows:

$$q(\mathbf{f}_{t-1}|\mathbf{f}_t, \mathbf{f}_x) = \mathcal{N}(\mathbf{f}_{t-1}; \boldsymbol{\mu}_{t \rightarrow t-1}, \boldsymbol{\sigma}_{t \rightarrow t-1}^2\mathbf{I}), \quad (15)$$

where:

$$\begin{aligned} q(\mathbf{f}_{t-1}|\mathbf{f}_t, \mathbf{f}_x) &\propto q(\mathbf{f}_{t-1}|\mathbf{f}_x)q(\mathbf{f}_t|\mathbf{f}_{t-1}, \mathbf{f}_x) = \\ q(\mathbf{f}_{t-1}|\mathbf{f}_x)q(\mathbf{f}_t|\mathbf{f}_{t-1}) &\Rightarrow \sigma_{t \rightarrow t-1} = \sigma_{t|t-1}\sigma_{t-1}/\sigma_t, \\ \boldsymbol{\mu}_{t \rightarrow t-1} &= (\alpha_{t|t-1}\sigma_{t-1}^2/\sigma_t^2)\mathbf{f}_t + (\alpha_{t-1}\sigma_{t|t-1}^2/\sigma_t^2)\mathbf{f}_x \end{aligned} \quad (16)$$

As discussed, the true denoising process can be approximated using a learned denoising distribution,  $p_\theta(\mathbf{f}_{t-1}|\mathbf{f}_t)$ .

## 3. Algorithms

Algorithms 1 and 2 summarize the training and decoding procedures of our neural codec.

## 4. Architecture of Diffusion-based Decoder

Fig. 1 illustrates our diffusion-based decoder design, employing a U-Net architecture for the diffusion model [1], incorporating ResNet blocks and self-attention modules. We've employed six units for both encoding and decoding within the U-Net framework. In the encoding pathway, the channel dimension is determined from the set

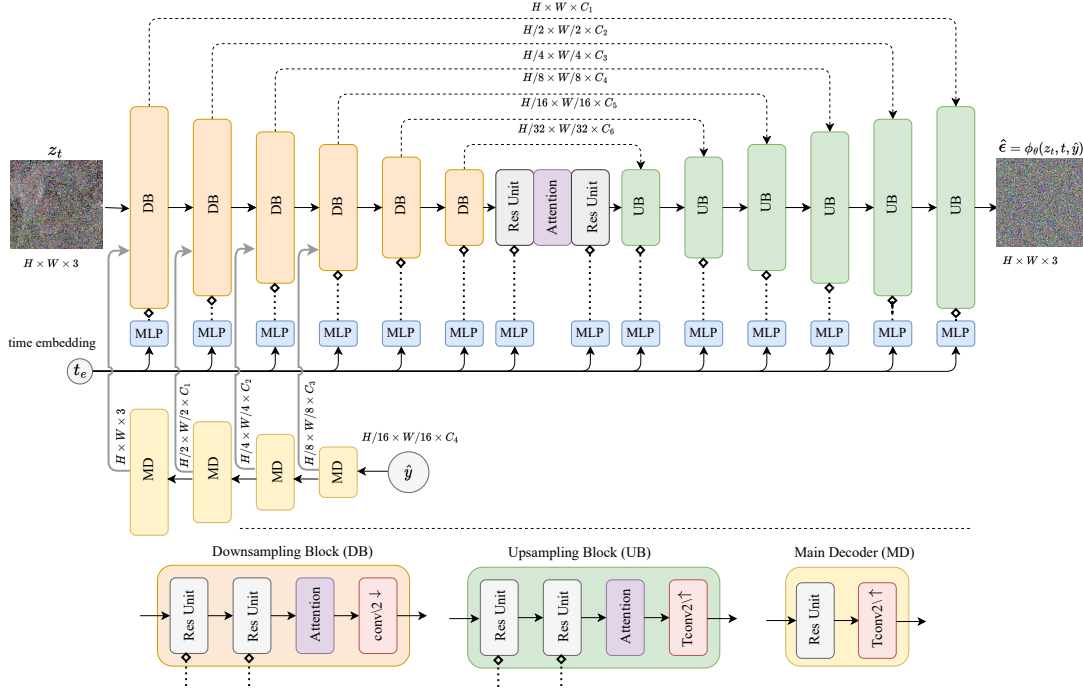


Figure 1. Architect of diffusion-based decoder.  $W$  and  $H$  correspond to the width and height of the input image, respectively.

### Algorithm 1 Training Neural Codec

Sample  $x \sim \text{dataset}$   
**repeat**  
 Sample  $t \sim \mathcal{U}(0, T)$   
 Sample  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
 $z_t = \mathbf{V} \alpha_t \mathbf{V}^T x + \mathbf{V} \sigma_t \mathbf{V}^T \epsilon$   
 $\tilde{y} \sim \mathcal{U}(\text{En}_\zeta(x) - 0.5, \text{En}_\zeta(x) + 0.5)$   
 $\hat{x}_t = \mathbf{V}(1/\alpha_t)(\mathbf{V}^T z_t - \sigma_t \mathbf{V}^T \phi_\theta(z_t, t, \tilde{y}))$   
 $L_{Dif} = \|\epsilon - \phi_\theta(z_t, t, \tilde{y})\|^2$   
 $L_T = (1 - \beta)L_{Dif} + \beta d_{\text{LPIPS}}(x, \hat{x}_t) - \lambda \log p_{\tilde{y}}(\tilde{y})$   
 $(\zeta, \theta) = (\zeta, \theta) - \eta \nabla_{\zeta, \theta} L_T$  ( $\eta$ : Learning Rate)  
**until** converged

### Algorithm 2 Decoding Compressed File

$\hat{y} \leftarrow$  Entropy decoded binary file using entropy model  
 $p_{\hat{y}}(\hat{y})$   
 Sample  $z_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
**for**  $t = T, \dots, 1$  **do**  
 $f_t = \mathbf{V}^T z_t$  and  $f_{\hat{\epsilon}} = \mathbf{V}^T \phi_\theta(z_t, t, \hat{y})$   
 $\sigma_{t \rightarrow t-1} = \sigma_{t|t-1} \sigma_{t-1} / \sigma_t$   
 $\hat{\mu}_{t \rightarrow t-1} = \frac{\alpha_{t|t-1} \sigma_{t-1}^2}{\sigma_t^2} f_t + \frac{\sigma_{t|t-1}^2}{\alpha_{t|t-1} \sigma_t^2} (f_t - \sigma_t f_{\hat{\epsilon}})$   
 $z_{t-1} \leftarrow \mathbf{V}(\hat{\mu}_{t \rightarrow t-1} + \sigma_{t \rightarrow t-1} f_{\hat{\epsilon}})$   
**end for**  
**Return**  $\hat{x} = z_0$

$\{C_1 = 64, C_2 = 128, C_3 = 192, C_4 = 256, C_5 = 320, C_6 = 384\}$ . The decoding process mirrors the encoding process in reverse. The main decoder (MD) comprises ResNet blocks and transposed convolutions, which serve to upscale the quantized latent representation  $\hat{y}$  to match the spatial dimensions of the inputs from the initial 4 U-Net encoding units. This setup enables us to introduce conditioning by concatenating the output of the main decoder layers with the input from the corresponding U-Net layer.

The time step  $t$  is initially linearly embedded into a vector with a dimension of 64. Subsequently, the resulting time embedding  $t_e$  is further processed through MLP layers, which are responsible for expanding it to align with the

channel size of the corresponding DB/UB layers.

## 5. Additional Qualitative Comparisons

As shown in Fig. 2, our model tends to generate fewer artifacts and is capable of decoding images with greater realism compared to both the HiFiC [3] and CDC [5] networks, even when using a significantly lower bit-rate.

## References

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2
- [2] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan



Figure 2. Additional Visual comparison of our method to the HiFiC and CDC models.

Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021. [1](#)

[3] Fabian Mentzer, George D Toderici, Michael Tschannen, and Eirikur Agustsson. High-fidelity generative image compression. *Advances in Neural Information Processing Systems*, 33:11913–11924, 2020. [3](#)

[4] Severi Rissanen, Markus Heinonen, and Arno Solin. Generative modelling with inverse heat dissipation. *arXiv preprint arXiv:2206.13397*, 2022. [2](#)

[5] Ruihan Yang and Stephan Mandt. Lossy image compression with conditional diffusion models. *arXiv preprint*

*arXiv:2209.06950*, 2022. [3](#)