

# IMPROVING VECTOR-QUANTIZED IMAGE MODELING WITH LATENT CONSISTENCY-MATCHING DIFFUSION

Bac Nguyen<sup>1\*</sup>, Chieh-Hsin Lai<sup>1</sup>, Yuhta Takida<sup>1</sup>, Naoki Murata<sup>1</sup>, Toshimitsu Uesaka<sup>1</sup>, Stefano Ermon<sup>2</sup>, Yuki Mitsufuji<sup>1,3</sup>

<sup>1</sup>Sony AI, <sup>2</sup>Stanford University, <sup>3</sup>Sony Group Corporation

## ABSTRACT

By embedding discrete representations into a continuous latent space, we can leverage continuous-space latent diffusion models to handle generative modeling of discrete data. However, despite their initial success, most latent diffusion methods rely on fixed pretrained embeddings, limiting the benefits of joint training with the diffusion model. While jointly learning the embedding (via reconstruction loss) and the latent diffusion model (via score matching loss) could enhance performance, end-to-end training risks embedding collapse, degrading generation quality. To mitigate this issue, we introduce VQ-LCMD, a continuous-space latent diffusion framework within the embedding space that stabilizes training. VQ-LCMD uses a novel training objective combining the joint embedding-diffusion variational lower bound with a consistency-matching (CM) loss, alongside a shifted cosine noise schedule and random dropping strategy. Experiments on several benchmarks show that the proposed VQ-LCMD yields superior results on FFHQ, LSUN Churches, and LSUN Bedrooms compared to discrete-state latent diffusion models. In particular, VQ-LCMD achieves an FID of 6.81 for class-conditional image generation on ImageNet with 50 steps.

## 1 INTRODUCTION

Vector-quantized variational autoencoders (VQ-VAE) (Van Den Oord et al., 2017; Razavi et al., 2019) have proven the usefulness of discrete latent representations in image generation (Gu et al., 2022; Chang et al., 2022). It typically involves training an encoder that compresses the image into a low-dimensional discrete latent space and then using a generative model such as autoregressive models (ARs) (Bengio et al., 2000; Brown et al., 2020) to learn and sample from this discrete latent space. Although ARs appear to dominate discrete data modeling, generating samples from these models incurs significant computational costs. Moreover, controllability is often challenging because the generation order has to be predetermined (Lou et al., 2023), making them less suitable for control tasks such as infilling and inpainting.

On the other hand, continuous-state diffusion models (CSDMs) (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020b) have shown promise as they enable efficient and rapid sampling without relying on the sequential attention mechanism of ARs. Diffusion models learn the inverse of a Markov chain that gradually converts data into pure Gaussian noise, using noise-conditioned score functions (i.e., gradients of log density), which are defined only for continuous data. The core concept is to progressively recover the original data by reversing the diffusion process. Diffusion models are notable for their high-fidelity generation (Dhariwal & Nichol, 2021; Lai et al., 2023a;b). They offer stable and relatively efficient training procedures that contribute to their success. Recent advances, such as consistency models (Song et al., 2023; Kim et al., 2023; Luo et al., 2023), have further enhanced diffusion models by reducing the number of sampling steps, making them more practical for real-world applications.

Despite the widespread popularity of CSDMs, their extension to discrete data remains limited. Previous attempts to address this limitation (Austin et al., 2021; Hoogeboom et al., 2021; Campbell et al., 2022; Sun et al., 2023; Lou et al., 2023) have focused on discrete-state diffusion models

\*Correspondence to bac.nguyencong@sony.com

(DSDMs), which define discrete corruption processes for discrete data and mimic Gaussian kernels used in continuous space. For instance, D3PMs (Austin et al., 2021) implemented the corruption process as random masking or token swapping and learned to reverse this process from the noisy data. However, unlike continuous diffusion processes, these corruption techniques do not progressively erase the semantic meaning of the data, potentially complicating the learning of the reverse procedure.

Alternatively, discrete data can be mapped into a continuous embedding space (Vahdat et al., 2021; Rombach et al., 2022; Sinha et al., 2021), followed by the application of CSDMs with typical Gaussian kernels, which enables progressive learning signals (Ho et al., 2020) and fine-grained sampling. This approach has been successful in various domains. However, it may not inherently yield satisfactory results (Li et al., 2022; Strudel et al., 2022; Dieleman et al., 2022). First, it requires a well-trained embedding for each new discrete dataset (Li et al., 2022) before training CSDMs. Since the embedding space and the denoising model are not trained end-to-end, this can result in suboptimal performance. Second, jointly training both components is challenging and prone to the embedding collapse problem (Dieleman et al., 2022; Gao et al., 2024), where all embeddings converge to similar vectors. While this convergence helps the diffusion model predict clean embeddings, it does not result in a meaningful model and instead leads to poor generation. To alleviate embedding collapse, previous work have explored normalizing embedding vectors to a fixed bounded norm (Dieleman et al., 2022) or mapping the predicted embedding to its nearest neighbor within the finite set of vectors (Li et al., 2022). However, the aforementioned manipulations may not yield satisfactory results in practice.

In response, this paper presents **Vector-Quantized Latent Consistency-Matching Diffusion (VQ-LCMD)**, a model that enables training of CSDM in discrete vector quantized latent space. We first compress images with a VQ-VAE into discrete tokens and then apply continuous diffusion to their embeddings. Our key contributions are summarized as follows.

- (i) A novel training objective is proposed to stabilize joint training of the embedding and diffusion variational lower bound. In particular, we enforce a *consistency-matching (CM)* loss that requires the model predictions to remain consistent over time. This ensures that the model produces stable outputs throughout the generation process, thereby helping to stabilize training.
- (ii) We identify several effective techniques to further enhance the generation quality. Specifically, we adopt (1) a shifted cosine noise schedule and (2) random embedding dropout. In addition, we perform a comprehensive analysis to evaluate the empirical impact of these techniques.
- (iii) Experiments on both unconditional and conditional image generation benchmarks are conducted to evaluate VQ-LCMD. The results show that VQ-LCMD effectively mitigates the embedding collapse issue and outperforms several baseline methods. VQ-LCMD achieves FID scores of 7.25 on FFHQ, 4.99 on LSUN Churches, 4.16 on LSUN Bedrooms, and 6.81 on ImageNet  $256 \times 256$ .

## 2 RELATED WORK

**Discrete-State Diffusion Models (DSDMs).** The idea is to establish a similar iterative refinement process for discrete data. The corruption process involves transitioning discrete values from one to another. This concept was initially introduced by Sohl-Dickstein et al. (2015) for binary sequence problems. Later, it was extended in multinomial diffusion (Hoogeboom et al., 2021). Austin et al. (2021) improved discrete diffusion by introducing diverse corruption processes, going beyond uniform transition. Based on the former framework, several extensions have been introduced for image modeling, e.g., MaskGIT (Chang et al., 2022), VQ-Diffusion (Gu et al., 2022), Token-Critic (Lezama et al., 2022), Muse (Chang et al., 2023), and Paella (Rampas et al., 2022). Additionally, Campbell et al. (2022) utilized Continuous Time Markov Chains for discrete diffusion. Despite their initial success, the corruptions introduced by these methods are characterized by their coarse-grained nature, making them inadequate for effectively modeling the semantic correlations between tokens.

**Continuous-Space Diffusion Models (CSDMs).** Li et al. (2022) addressed the challenge of controlling language models with Diffusion-LM, a non-autoregressive language model based on continuous diffusion. A similar idea has been introduced in SED (Strudel et al., 2022), DiNoiSer (Ye et al., 2023), CDCD (Dieleman et al., 2022), Bit Diffusion (Chen et al., 2022), Plaid (Gulrajani & Hashimoto, 2024), and Difformer (Gao et al., 2024). The challenge of end-to-end training for both embeddings

and CSDMs has not been fully addressed in these methods. To avoid embedding collapse, existing techniques either normalize the embeddings (Dieleman et al., 2022) or use heuristic methods (Li et al., 2022), which are not generally effective and may lead to training instability (Dieleman et al., 2022; Strudel et al., 2022). Recently, Lou et al. (2023) and Meng et al. (2022) proposed an alternative concrete score function for discrete settings, which captures the surrogate “gradient” information within discrete spaces.

### 3 PRELIMINARY

In image synthesis, directly modeling raw pixels can be computationally expensive, especially for high-resolution images. To reduce this cost, the training process can be divided into two phases. First, an autoencoder is trained to produce lower-dimensional representations, followed by training a generative model in this latent space. This is because pixel-based representations of images contain high-frequency details but little semantic variation (Rombach et al., 2022). Below, we outline the concept of VQ-VAE and reformulate diffusion models within this discrete latent space.

#### 3.1 DISCRETE REPRESENTATION OF IMAGES

To compress an image into discrete representations, VQ-VAE (Van Den Oord et al., 2017) employs a learnable discrete codebook combined with nearest neighbor search to train an encoder-decoder architecture. The nearest neighbor search is performed between the encoder output and the latent embeddings in the codebook. Finally, the resulting discrete latent sequence is then passed to the decoder to reconstruct the image. To further improve the generation fidelity, VQGAN (Esser et al., 2021) leverages adversarial training to the decoder output.

Given an image, we obtain a sequence of discrete image tokens  $\mathbf{x} = [x_1, \dots, x_M]$  with a pre-trained VQ-VAE, where each image token  $x_i$  belongs to one of the  $K$  categories  $\{1, \dots, K\}$  in the codebook. Here  $M$  denotes the number of image tokens in the discrete space. The distribution over discrete latent variables  $\mathbf{x}$  is multinomial and denoted as  $P(\mathbf{x})$ .

#### 3.2 DIFFUSION MODELS IN DISCRETE SPACE

Our goal is to learn a generative model that approximates the probability mass function  $P(\mathbf{x})$ . To handle discontinuity, we propose using continuous embeddings, where different categories are represented by real-valued vectors. Let  $\phi = \{e_1, \dots, e_K\}$ , where  $e_k \in \mathbb{R}^D$ , be a set of vectors, the embeddings of  $\mathbf{x}$  are then defined as  $\Psi_\phi(\mathbf{x}) = [e_{x_1}, \dots, e_{x_M}]$ . We define a sequence of increasingly noisy versions of  $\Psi_\phi(\mathbf{x})$  as  $\mathbf{z}_t$ , where  $t$  ranges from  $t = 0$  (least noisy) to  $t = 1$  (most noisy). Next, we adopt the variational diffusion formulation (Kingma et al., 2021), incorporating the embedding  $\Psi_\phi$ .

**Forward process.** We define the forward process as a Markov chain, which progressively corrupts the data with Gaussian noise (Kingma et al., 2021; Ho et al., 2020). For any  $t \in [0, 1]$ , the conditional distribution of  $\mathbf{z}_t$  given  $\mathbf{x}$  is modeled as

$$q_\phi(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\mathbf{z}_t|\alpha_t\Psi_\phi(\mathbf{x}), \sigma_t^2\mathbf{I}),$$

where  $\alpha_t$  and  $\sigma_t$  are positive scalar-value functions of  $t$ , which determine how much noise is added to the embeddings of  $\mathbf{x}$ . We consider a variance-preserving process, i.e.,  $\alpha_t^2 + \sigma_t^2 = 1$ . The marginal distribution  $q_\phi(\mathbf{z}_t)$  is a mixture of Gaussian distributions. Due to the Markovian property by construction, the transition probability distributions are given by

$$q(\mathbf{z}_t|\mathbf{z}_s) = \mathcal{N}(\mathbf{z}_t|\alpha_{t|s}\mathbf{z}_s, \sigma_{t|s}^2\mathbf{I}),$$

where  $\alpha_{t|s} = \alpha_t/\alpha_s$  and  $\sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2\sigma_s^2$ . Conditioned on the clean discrete variable  $\mathbf{x}$ , the forward process posterior distribution is derived as

$$q_\phi(\mathbf{z}_s|\mathbf{z}_t, \mathbf{x}) = \mathcal{N}(\mathbf{z}_s|\mu_\phi(\mathbf{z}_t, \mathbf{x}; s, t), \sigma^2(s, t)\mathbf{I}),$$

where  $\mu_\phi(\mathbf{z}_t, \mathbf{x}; s, t) = \frac{\alpha_{t|s}\sigma_s^2}{\sigma_t^2}\mathbf{z}_t + \frac{\alpha_s\sigma_{t|s}^2}{\sigma_t^2}\Psi_\phi(\mathbf{x})$  and  $\sigma^2(s, t) = \sigma_{t|s}^2\sigma_s^2/\sigma_t^2$ .

**Reverse process.** We gradually denoise the latent variables toward the data distribution by a Markov process. Starting from the standard Gaussian prior  $p(\mathbf{z}_1)$ , the Markov reverse process runs backward

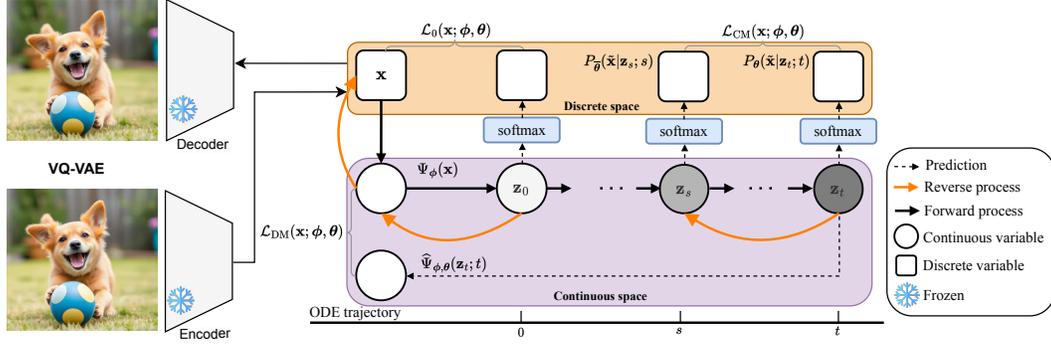


Figure 1: Training procedure of VQ-LCMD. An image is compressed into a sequence of discrete tokens  $\mathbf{x}$  using a pre-trained VQ-VAE. VQ-LCMD learns to generate the discrete latent representations  $\mathbf{x}$  using the consistency-matching (CM) loss, diffusion loss, and reconstruction loss.

from  $t = 1$  to  $t = 0$ . Let  $\theta$  denote the parameters of the denoising model, the conditional probability distribution  $p_{\phi, \theta}(\mathbf{z}_s | \mathbf{z}_t; s, t)$  for any  $0 \leq s \leq t \leq 1$  in the reverse diffusion process is parameterized by a Gaussian. More specifically, it is given by

$$p_{\phi, \theta}(\mathbf{z}_s | \mathbf{z}_t; s, t) = \mathcal{N}(\mathbf{z}_s | \hat{\mu}_{\phi, \theta}(\mathbf{z}_t; s, t), \sigma^2(s, t) \mathbf{I}), \quad (1)$$

where  $\hat{\mu}_{\phi, \theta}(\mathbf{z}_t; s, t) = \frac{\alpha_t |s| \sigma_s^2}{\sigma_t^2} \mathbf{z}_t + \frac{\alpha_s \sigma_t^2 |s|}{\sigma_t^2} \hat{\Psi}_{\phi, \theta}(\mathbf{z}_t; t)$  and  $\hat{\Psi}_{\phi, \theta}(\mathbf{z}_t; t)$  denotes the predicted embeddings of  $\Psi_{\phi}(\mathbf{x})$  based on its noisy version  $\mathbf{z}_t$ .

**Network parametrization.** We parameterize  $\hat{\Psi}_{\phi, \theta}(\mathbf{z}_t; t)$  as an average over embeddings, where the  $i$ -element of  $\hat{\Psi}_{\phi, \theta}(\mathbf{z}_t; t)$  is given by

$$[\hat{\Psi}_{\phi, \theta}(\mathbf{z}_t; t)]_i = \sum_{k=1}^K P_{\theta}(\tilde{x}_i = k | \mathbf{z}_t; t) \mathbf{e}_k.$$

As the forward process factorizes across  $M$  tokens  $q_{\phi}(\mathbf{z}_s | \mathbf{z}_t, \mathbf{x}) = \prod_{i=1}^M q_{\phi}(\mathbf{z}_{s,i} | \mathbf{z}_{t,i}, x_i)$ , we also model the reverse process as a factorized distribution. In particular, to estimate the posterior probability  $P_{\theta}(\tilde{\mathbf{x}} | \mathbf{z}_t; t)$ , we use a neural network  $f_{\theta}(\mathbf{z}_t; t)$  to predict  $K$  logits for each token, followed by a softmax nonlinearity, i.e.,

$$P_{\theta}(\tilde{\mathbf{x}} | \mathbf{z}_t; t) = \prod_{i=1}^M \text{softmax}([f_{\theta}(\mathbf{z}_t; t)]_i).$$

As  $t$  approaches zero, the decoding process from  $\mathbf{z}_0$  to  $\mathbf{x}$  gives a learning signal for  $\phi$ .

**Variational lower bound.** Following (Kingma et al., 2021), the negative variational lower bound (VLB) is derived as

$$-\log P_{\phi, \theta}(\mathbf{x}) \leq \mathbb{E}_{\epsilon} [-\log P_{\theta}(\mathbf{x} | \mathbf{z}_0; 0)] + D_{\text{KL}}(q_{\phi}(\mathbf{z}_1 | \mathbf{x}) || p(\mathbf{z}_1)) + \mathcal{L}_{\infty}(\mathbf{x}; \phi, \theta), \quad (2)$$

where  $\mathbf{z}_t = \alpha_t \Psi_{\phi}(\mathbf{x}) + \sigma_t \epsilon$  with  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and the diffusion loss is simplified to

$$\mathcal{L}_{\infty}(\mathbf{x}; \phi, \theta) = -\frac{1}{2} \mathbb{E}_{\epsilon, t} [\text{SNR}(t)' \|\Psi_{\phi}(\mathbf{x}) - \hat{\Psi}_{\phi, \theta}(\mathbf{z}_t; t)\|^2]$$

with  $\text{SNR}(t) = \alpha_t^2 / \sigma_t^2$  the signal-to-noise ratio. Under certain conditions<sup>1</sup>, the prior loss is close to zero as  $q_{\phi}(\mathbf{z}_1 | \mathbf{x}) \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Unlike CSDMs, the reconstruction loss in our case  $\mathcal{L}_0(\mathbf{x}; \phi, \theta) = \mathbb{E}_{\epsilon} [-\log P_{\theta}(\mathbf{x} | \mathbf{z}_0; 0)]$  is important since it involves both denoising and embedding parameters. A remarkable result shown by Kingma et al. (2021) is that the diffusion loss is invariant to the noise schedule except at  $t = 0$  and  $t = 1$ .

Although  $\phi$  and  $\theta$  can be jointly trained by minimizing Eq. (2), this approach often results in a solution in which most embeddings collapse into nearly identical vectors with minimal variance, leading to degraded generation quality (refer to our ablation studies in Table 1).

<sup>1</sup>In theory, we require that  $\alpha_1 \Psi_{\phi}(\mathbf{x}) = \mathbf{0}$  to ensure that the prior loss is equal zero.

## 4 PROPOSED METHOD

This section presents VQ-LCMD. First, we introduce the consistency-matching (CM) loss to ensure consistent probability predictions across timesteps. Next, we propose re-weighting the objectives in the loss function, along with an improved noise schedule and a random dropping strategy, to further improve results. The overall training and objective function of VQ-LCMD is illustrated in Fig. 1.

### 4.1 CONSISTENCY-MATCHING LOSS

Considering  $\Psi_\phi(\mathbf{x})$  as clean data in the continuous space, the evolution of  $\Psi_\phi(\mathbf{x})$  over time can be described by the probability flow ordinary differential equation (PF-ODE) (Song et al., 2020b). This PF-ODE allows a deterministic bijection between the embeddings  $\Psi_\phi(\mathbf{x})$  and latent representations  $\mathbf{z}_t$ . Intuitively, a random noise perturbation  $\mathbf{z}_t$  of  $\Psi_\phi(\mathbf{x})$  and its relatively nearby point  $\bar{\mathbf{z}}_s$  along the same trajectory should yield nearly the same prediction. To ensure these consistent outputs for arbitrary  $\mathbf{z}_t$ , we propose the consistency-matching (CM) loss

$$\mathcal{L}_{\text{CM}}(\mathbf{x}; \phi, \theta) = \mathbb{E}_{\epsilon, t, s} [D_{\text{KL}}(P_{\bar{\theta}}(\tilde{\mathbf{x}}|\bar{\mathbf{z}}_s; s) \| P_{\theta}(\tilde{\mathbf{x}}|\mathbf{z}_t; t))], \quad (3)$$

where  $\bar{\theta}$  denotes an exponential moving average (EMA), i.e.,  $\bar{\theta} \leftarrow \text{stopgrad}(\eta\bar{\theta} + (1-\eta)\theta)$  with a decay rate of  $\eta$ . The time variables are sampled uniformly, where  $t \sim \mathcal{U}(0, 1)$  and  $s$  is sampled from the interval  $[0, t]$ , i.e.,  $s \sim \mathcal{U}(0, t)$ . Here,  $\mathbf{z}_t$  is obtained by perturbing  $\Psi_\phi(\mathbf{x})$  to the noise level  $t$  using the transition kernel  $q_\phi(\mathbf{z}_t|\mathbf{x})$  and  $\bar{\mathbf{z}}_s$  is obtained by taking a PF-ODE step using the EMA model. There are many several ODE solvers to get the PF-ODE step such as Euler (Song et al., 2020b) and Heun (Karras et al., 2022) solvers. For simplicity, we use the DDIM sampler (Song et al., 2020a), which applies the Euler discretization on the PF-ODE. Under variance preserving settings, it is computed as

$$\bar{\mathbf{z}}_s = \alpha_s \Psi_{\bar{\phi}}(\mathbf{x}) + (\sigma_s/\sigma_t)(\mathbf{z}_t - \alpha_t \Psi_{\bar{\phi}}(\mathbf{x})),$$

where  $\bar{\phi} \leftarrow \text{stopgrad}(\eta\bar{\phi} + (1-\eta)\phi)$ .

Our intuition for the CM loss is that when the timesteps  $t$  are small, the model learns the true distribution through the reconstruction loss. As training progresses, this consistency is propagated to later timesteps, eventually reaching  $t = 1$ . The CM loss encourages the probability distributions in neighboring latent variables to converge. Once the model is fully trained, it consistently produces the same probability distribution across the entire PF-ODE trajectory. Since the reconstruction loss enforces the mapping from the embedding space back to discrete data, the learning signal from the reconstruction loss is propagated through the entire PF-ODE trajectory.

**Connection of CM loss to existing works.** When the distribution  $P(\mathbf{x})$  is continuous, Eq. (3) recovers the consistency training objective in CSDMs (Song et al., 2023; Kim et al., 2023; Lai et al., 2023b), which matches clean predictions from models along the same sampling PF-ODE trajectory. Specifically, for any noisy sample  $\mathbf{z}_t$  at time  $t$ ,  $P_{\theta}(\tilde{\mathbf{x}}|\mathbf{z}_t; t)$  serves as a deterministic consistency function (Song et al., 2023)  $\mathbf{h}_{\theta}(\mathbf{z}_t; t)$  predicting the clean sample at time 0 from  $\mathbf{z}_t$ , regarded as a normal distribution centered around  $\mathbf{h}_{\theta}(\mathbf{z}_t; t)$  with small variance. Thus, using the closed-form KL divergence of two normal distributions, Eq. (3) becomes:

$$\mathcal{L}_{\text{CM}}(\mathbf{x}; \phi, \theta) \propto \mathbb{E}_{\epsilon, t, s} \left[ \left\| \mathbf{h}_{\bar{\theta}}(\mathbf{z}_s; s) - \mathbf{h}_{\theta}(\mathbf{z}_t; t) \right\|_2^2 \right],$$

which coincides with a special case of "soft consistency" proposed by Kim et al. (2023) (with their intermediate timesteps  $u$  and end time  $s$  replaced by our alternate starting time  $s$ , and our end time 0). Here,  $\propto$  denotes the omission of multiplicative or additive constants that are independent of the training parameters.

### 4.2 FINAL LOSS FUNCTION

Although  $-\text{SNR}(t)'$  in Eq. (2) provides the correct scaling to treat the objective function as an upper bound of the negative log-likelihood, we hypothesize this weighting function may disrupt the balance between training the reconstruction loss and diffusion loss in practice. Instead of minimizing directly the diffusion loss, we simplify it as

$$\mathcal{L}_{\text{DM}}(\mathbf{x}; \phi, \theta) = \mathbb{E}_{\epsilon, t} \left[ \left\| \Psi_\phi(\mathbf{x}) - \hat{\Psi}_{\phi, \theta}(\mathbf{z}_t; t) \right\|_2^2 \right].$$

This ensures that the loss is evenly distributed over different timesteps. The rationale is that alleviating the error in a large noise level can help the model avoid constant embeddings (Li et al., 2022).

Putting it all together, the overall objective function of VQ-LCMD is given by

$$\min_{\phi, \theta} \mathbb{E}_{\mathbf{x}}[\mathcal{L}(\mathbf{x}; \phi, \theta)] = \mathbb{E}_{\mathbf{x}}[\mathcal{L}_0(\mathbf{x}; \phi, \theta) + \beta_{\text{DM}}\mathcal{L}_{\text{DM}}(\mathbf{x}; \phi, \theta) + \beta_{\text{CM}}\mathcal{L}_{\text{CM}}(\mathbf{x}; \phi, \theta)],$$

where  $\beta_{\text{DM}} \geq 0$  and  $\beta_{\text{CM}} \geq 0$  are hyperparameters. By tuning  $\beta_{\text{DM}}$  and  $\beta_{\text{CM}}$ , we can find the right balance between the objective functions.

### 4.3 NOISE SCHEDULE

Although the diffusion loss remains invariant to the noise schedule (Kingma et al., 2021), it is essential to determine how noise evolves during the diffusion process (Song et al., 2021; Kingma & Gao, 2023). This is because Monte Carlo sampling is employed to estimate the diffusion loss, and thus the training dynamics are influenced by the choice of noise schedule. If the embedding norms are large, denoising would be a trivial task for low noise levels. This is not desired because the denoising model has only a small time window to generate the global structure of the meaningful embedding. To address this, we use the shifted cosine noise schedule (Hoogeboom et al., 2023),

$$\log \text{SNR}(t) = -2 \log \tan(\pi t/2) + s,$$

where  $s \in \mathbb{R}$  is a hyper-parameter. This adjustment changes the noise schedule by shifting its log SNR curve. In particular, when  $s = 0$ , it corresponds to the cosine noise schedule (Nichol & Dhariwal, 2021). Essentially, the noise schedule implies different weights in the diffusion loss per noise level (Kingma & Gao, 2023). As illustrated in Fig. 2, by moving the curve to the left, it gives more importance for higher degrees of noise.

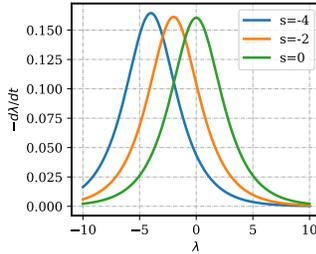


Figure 2: Shifted cosine noise schedule with different shifting factors  $s$ ,  $\lambda(t) = \log \text{SNR}(t)$ .

### 4.4 RANDOM DROPPING

Given the noised embeddings,  $\mathcal{L}_{\text{CM}}$  used in our training objective ensures the same prediction for the posterior probability  $P_{\theta}(\tilde{\mathbf{x}}|\mathbf{z}_t; t)$  at any timestep. During joint training, it encourages the model to distinguish the embeddings by increasing their parameter magnitudes. To avoid this shortcut solution, we propose to randomly drop the embeddings. This forces the representations to be more semantic (He et al., 2022). Let  $\mathbf{m}_{\text{RD}} \in \{0, 1\}^M$  denote a binary mask that indicates which tokens are replaced with a special [mask] token. During training, embeddings of  $\mathbf{x}$  become  $\Psi_{\phi}(\mathbf{x} \odot \mathbf{m}_{\text{RD}})$ . This is because only a portion of the embeddings is used to predict the other tokens. It requires the model to understand the relationship between masked and unmasked tokens. When similar tokens frequently appear in similar contexts, the model learns to associate these tokens closely in the embedding space, as their contextual meanings are similar.

## 5 EXPERIMENTS

This section evaluates the performance of VQ-LCMD on several benchmarks. We begin by outlining the experimental setups, followed by comprehensive experiments covering both conditional and unconditional image generation tasks. Finally, we provide detailed ablation studies to analyze VQ-LCMD.

### 5.1 EXPERIMENTAL SETUP

We briefly describe the datasets, baselines, and metrics used for evaluation. Additional details are provided in Appendix A.

**Datasets.** For unconditional generation, our benchmark consists of three datasets: FFHQ (Karras et al., 2019), LSUN Bedrooms, and LSUN Churches (Yu et al., 2015). The FFHQ dataset contains 70K examples of human faces, while the LSUN Bedrooms dataset contains 3M images of bedrooms,

Table 1: Results of ablation studies on the FFHQ dataset

Method	FID↓	Prec.↑	Rec.↑
CSDM	77.09	0.41	0.07
CSDM w $\ell_2$ -norm	52.35	0.55	0.12
VQ-LCMD w/o $\mathcal{L}_{CM}$	186.95	0.02	0.00
VQ-LCMD w/o NS	11.19	0.71	0.42
VQ-LCMD w/o RD	<u>8.20</u>	<b>0.73</b>	0.42
VQ-LCMD	<b>7.25</b>	<u>0.72</u>	<b>0.46</b>



Figure 3: VQ-LCMD samples for unconditional generation

Table 2: Results for unconditional generation on FFHQ, LSUN Churches, and LSUN Bedrooms. The scores of FID, Precision, and Recall are shown. The **best** and second best results are marked.

Method	FFHQ			LSUN Churches			LSUN Bedrooms		
	FID↓	Prec.↑	Rec.↑	FID↓	Prec.↑	Rec.↑	FID↓	Prec.↑	Rec.↑
<i>Discrete-Space Diffusion Models</i>									
D3PM Uniform	9.49	0.71	0.41	6.02	0.68	0.39	6.60	0.60	0.35
VQ-Diffusion	<u>8.79</u>	0.70	<u>0.43</u>	6.88	0.72	0.37	7.19	0.54	0.37
MaskGIT	11.45	<b>0.75</b>	0.42	<u>5.59</u>	0.65	<b>0.44</b>	8.39	0.66	0.33
<i>Continuous-Space Diffusion Models</i>									
CSDM <sup>†</sup>	12.66	<u>0.73</u>	0.38	7.88	<b>0.76</b>	0.36	<u>4.93</u>	<u>0.71</u>	<u>0.38</u>
VQ-LCMD (ours)	<b>7.25</b>	<u>0.72</u>	<b>0.46</b>	<b>4.99</b>	<u>0.75</u>	<u>0.42</u>	<b>4.16</b>	<b>0.72</b>	<b>0.40</b>

and the LSUN Church dataset contains 126K images of churches. For conditional generation, we use ImageNet (Deng et al., 2009). These datasets are widely used in the literature. All images have a resolution of  $256 \times 256$  and VQGAN (Esser et al., 2021) is used to downsample the images into discrete representations of  $16 \times 16$  with a codebook size of 1024.

**Baselines and metrics.** We evaluate VQ-LCMD against several baselines, including D3PM with uniform transition probabilities (Austin et al., 2021), VQ-Diffusion (Gu et al., 2022), and MaskGIT (Chang et al., 2022). Additionally, we include results for CSDM using fixed embeddings (CSDM<sup>†</sup>), where embeddings are initialized from the pretrained VQGAN codebook and remain fixed throughout training. For evaluation, we report the Fréchet Inception Distance (FID) between 50,000 generated images and real images. We also provide performance metrics in terms of Precision and Recall (Kynkäänniemi et al., 2019). For conditional image generation, we use the Inception Score (IS) as an additional metric to measure the image quality.

## 5.2 UNCONDITIONAL IMAGE GENERATION

Table 2 presents the results for unconditional image generation tasks. To make a fair comparison, all models are configured with 200 steps for inference. VQ-LCMD consistently achieves the lowest FID scores. Furthermore, we investigate the impact of using pretrained embeddings in CSDM<sup>†</sup> and demonstrate that while it yields satisfactory results, employing trainable embeddings significantly enhances the performance. On LSUN Bedrooms, VQ-LCMD outperforms the baseline methods by a substantial margin, achieving the highest Precision and Recall scores. These findings underline the superiority of VQ-LCMD in generating high-quality samples. The observed improvements in our method compared to discrete diffusion baselines confirm that continuous diffusion models can provide an effective solution for discrete data. Fig. 3 illustrates samples generated by VQ-LCMD.

## 5.3 CONDITIONAL IMAGE GENERATION

Table 3 presents the results for class-conditional image generation tasks. To improve the sample quality of conditional diffusion models, we employ the classifier-free guidance (Ho & Salimans, 2021). Essentially, it guides the sampling trajectories toward higher-density data regions. During training, we randomly drop 10% of the conditions and set the dropped conditions to the null token. Our method achieves a FID of 6.81 and an IS of 225.31 with 50 sampling steps. VQ-LCMD notably

outperforms both VQGAN and VQVAE-2 by a substantial margin. Compared to MaskGIT<sup>2</sup>, VQ-LCMD provides competitive FID results and exceeds in IS. However, it is important to note, as highlighted by Besnier and Chen (Besnier & Chen, 2023), that MaskGIT requires specific sampling adjustments, such as adding Gumbel noise with a linear decay, to improve its FID. In contrast, VQ-LCMD operates without such sampling heuristics. In addition, VQ-LCMD performs better than VQ-Diffusion in both FID and IS metrics. For reference samples generated by VQ-LCMD, please refer to Appendix E.

Table 3: Comparison with generative models on ImageNet  $256 \times 256$ . The results of the existing methods are obtained from their respective published works.

Model	# params	# steps	FID↓	IS↑	Precision↑	Recall↑
VQGAN (Esser et al., 2021)	1.4B	256	15.78	74.3	n/a	n/a
MaskGIT (Besnier & Chen, 2023) (PyTorch)	246M	8	<b>6.80</b>	<u>214.0</u>	0.82	0.51
VQVAE-2(Razavi et al., 2019)	13.5B	5120	31.11	45.00	0.36	<u>0.57</u>
BigGAN-deep (Brock et al., 2019)	160M	1	6.95	198.2	<b>0.87</b>	0.28
Improved DDPM (Nichol & Dhariwal, 2021)	280M	250	12.26	n/a	0.70	<b>0.62</b>
VQ-Diffusion (Gu et al., 2022)	518M	100	11.89	n/a	n/a	n/a
VQ-LCMD (ours)	246M	50	<u>6.81</u>	<b>225.31</b>	<u>0.84</u>	0.38

#### 5.4 ABLATION STUDIES

This section presents ablation studies. For additional analysis, please see Appendix C. We investigate the impact of individual components introduced in VQ-LCMD on overall performance. Specifically, we examine the shifted cosine noise schedule (NS), random dropping (DR), and consistency-matching loss ( $\mathcal{L}_{CM}$ ). The results are presented in Table 1. The baseline method CSDM, trained by minimizing Eq. (2), is unable to generate meaningful images. While incorporating an  $\ell_2$ -norm regularization on the embeddings provides some improvement, it does not completely resolve the collapse issue. VQ-LCMD (incorporating our novel components  $\mathcal{L}_{CM} + NS + RD$ ) achieves the best performance. Without RD, the model produces inferior results. Removing NS leads to notable performance degradation. On the other hand, omitting  $\mathcal{L}_{CM}$  results in embedding collapse. These findings highlight the essential role of each component in mitigating the embedding collapse and improving overall performance.

## 6 CONCLUSION

We have introduced VQ-LCMD, a continuous diffusion model tailored for modeling discrete vector-quantized latent distributions, which jointly learns the embeddings and the denoising model. VQ-LCMD uses a novel training objective combining the joint embedding-diffusion variational lower bound with a consistency-matching (CM) loss, alongside a shifted cosine noise schedule and random dropping strategy. Experimental results show that VQ-LCMD not only alleviates the embedding collapse problem, but also exceeds baseline discrete-state diffusion models.

**Limitations and future work.** In this work, VQ-LCMD is implemented using the Transformer architecture, but we emphasize that the architecture choice is orthogonal to the proposed framework and can be extended to other architectures. Although our main focus is image generation task, VQ-LCMD can be applied to any task involving discrete data. Future work will focus on applying it to additional data types, such as graphs and text. It is also interesting to explore more advanced sampling techniques to improve the overall generation quality of VQ-LCMD.

## REFERENCES

Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. In *NeurIPS*, pp. 17981–17993, 2021.

<sup>2</sup>Since VQ-LCMD is implemented in PyTorch, we also use the PyTorch implementation (Besnier & Chen, 2023) of MaskGIT to ensure a fair comparison.

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. In *NeurIPS*, 2016.
- Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. A neural probabilistic language model. In *NeurIPS*, 2000.
- Victor Besnier and Mickael Chen. A Pytorch reproduction of masked generative image transformer. *arXiv preprint arXiv:2310.14400*, 2023.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, 2019.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, pp. 1877–1901, 2020.
- Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and Arnaud Doucet. A continuous time framework for discrete denoising models. In *NeurIPS*, pp. 28266–28279, 2022.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. MaskGIT: Masked generative image transformer. In *CVPR*, pp. 11315–11325, 2022.
- Huiwen Chang, Han Zhang, Jarred Barber, Aaron Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Patrick Murphy, William T. Freeman, Michael Rubinstein, Yanzhen Li, and Dilip Krishnan. Muse: Text-to-image generation via masked generative transformers. In *ICML*, pp. 4055–4075, 2023.
- Ting Chen, Ruixiang Zhang, and Geoffrey E. Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. In *ICLR*, 2022.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, pp. 248–255, 2009.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In *NeurIPS*, pp. 8780–8794, 2021.
- Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin, Pierre H Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, et al. Continuous diffusion for categorical data. *arXiv preprint arXiv:2211.15089*, 2022.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, pp. 12873–12883, 2021.
- Zhujin Gao, Junliang Guo, Xu Tan, Yongxin Zhu, Fang Zhang, Jiang Bian, and Linli Xu. Diffuser: Empowering diffusion model on embedding space for text generation. *arXiv preprint arXiv:2212.09412v3*, 2024.
- Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, pp. 10696–10706, 2022.
- Ishaan Gulrajani and Tatsunori B Hashimoto. Likelihood-based diffusion language models. In *NeurIPS*, volume 36, 2024.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, pp. 16000–16009, 2022.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS Workshop*, 2021.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, pp. 6840–6851, 2020.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *ICLR*, 2020.

- Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. In *NeurIPS*, pp. 12454–12465, 2021.
- Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. Simple diffusion: end-to-end diffusion for high resolution images. In *ICML*, pp. 13213–13232, 2023.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pp. 4401–4410, 2019.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *NeurIPS*, 2022.
- Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Yutong He, Yuki Mitsufuji, and Stefano Ermon. Consistency trajectory models: Learning probability flow ODE trajectory of diffusion. In *NeurIPS Workshop*, 2023.
- Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. In *NeurIPS*, pp. 21696–21707, 2021.
- Diederik P Kingma and Ruiqi Gao. Understanding diffusion objectives as the elbo with simple data augmentation. In *NeurIPS*, 2023.
- Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *Advances in neural information processing systems*, 32, 2019.
- Chieh-Hsin Lai, Yuhta Takida, Naoki Murata, Toshimitsu Uesaka, Yuki Mitsufuji, and Stefano Ermon. FP-diffusion: Improving score-based diffusion models by enforcing the underlying score Fokker-Planck equation. In *ICML*, pp. 18365–18398, 2023a.
- Chieh-Hsin Lai, Yuhta Takida, Toshimitsu Uesaka, Naoki Murata, Yuki Mitsufuji, and Stefano Ermon. On the equivalence of consistency-type models: Consistency models, consistent diffusion models, and fokker-planck regularization. In *ICML Workshop*, 2023b.
- José Lezama, Huiwen Chang, Lu Jiang, and Irfan Essa. Improved masked image generation with token-critic. In *ECCV*, pp. 70–86, 2022.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-LM improves controllable text generation. In *NeurIPS*, pp. 4328–4343, 2022.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion language modeling by estimating the ratios of the data distribution. *arXiv preprint arXiv:2310.16834*, 2023.
- Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching: Generalized score matching for discrete data. In *NeurIPS*, pp. 34532–34545, 2022.
- Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*, pp. 8162–8171, 2021.
- Dominic Rampas, Pablo Pernias, Elea Zhong, and Marc Aubreville. Fast text-conditional discrete denoising on vector-quantized latent spaces. *arXiv preprint arXiv:2211.07292*, 2022.
- Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with VQ-VAE-2. In *NeurIPS*, 2019.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10684–10695, 2022.
- Abhishek Sinha, Jiaming Song, Chenlin Meng, and Stefano Ermon. D2c: Diffusion-decoding models for few-shot conditional generation. *NeurIPS*, 34:12533–12548, 2021.

- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pp. 2256–2265, 2015.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020a.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2020b.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. In *NeurIPS*, pp. 1415–1428, 2021.
- Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.
- Robin Strudel, Corentin Tallec, Florent Althé, Yilun Du, Yaroslav Ganin, Arthur Mensch, Will Grathwohl, Nikolay Savinov, Sander Dieleman, Laurent Sifre, et al. Self-conditioned embedding diffusion for text generation. *arXiv preprint arXiv:2211.04236*, 2022.
- Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time discrete diffusion models. In *ICLR*, 2023.
- Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *NeurIPS*, 34:11287–11302, 2021.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Mingxuan Wang. Dinoiser: Diffused conditional sequence learning by manipulating noises. *arXiv preprint arXiv:2302.10025*, 2023.
- Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, 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.

# Appendix

---

<b>A Implementation Details</b>	<b>12</b>
<b>B Latent Variable Classifier-Free Guidance</b>	<b>12</b>
<b>C Ablation Studies</b>	<b>13</b>
C.1 Pretrained vs. Learnable Embeddings . . . . .	13
C.2 Number of Sampling Steps . . . . .	13
C.3 Dropping Strategies . . . . .	13
C.4 Weighting Terms . . . . .	14
C.5 Embedding dimensionality . . . . .	14
C.6 Ablation Studies on ImageNet . . . . .	15
<b>D Pseudocode</b>	<b>15</b>
<b>E Additional Samples</b>	<b>15</b>

---

## A IMPLEMENTATION DETAILS

The prediction network  $f_{\theta}(\mathbf{z}_t; t)$  is a bidirectional Transformer (Vaswani et al., 2017). For unconditional generation, the network consists of 15 layers, 8 attention heads, and 512 embedding dimensions (a total of 56M parameters). We apply a dropout rate of 0.1 to the self-attention layers. All models are trained on 4 NVIDIA DGX H100 GPUs with a batch size of 128. We use sinusoidal positional embeddings. For conditional generation on ImageNet, we scale up the model to 24 layers, 16 attention heads, and 768 embedding dimensions (a total of 246M parameters). Following (Gu et al., 2022), the conditional class label is injected into the model using Adaptive Layer Normalization (Ba et al., 2016) (AdaLN), i.e.,  $\text{AdaLN}(\mathbf{h}, t) = (1 + \mathbf{a}_t)\text{LayerNorm}(\mathbf{h}) + \mathbf{b}_t$ , where  $\mathbf{h}$  denotes the activation,  $\mathbf{a}_t$  and  $\mathbf{b}_t$  are obtained from a linear projection of the class embedding. We do not use any sampling heuristics such as top- $k$  or nucleus sampling (Holtzman et al., 2020). For random dropping, the dropping probability is fixed to 0.2 as the default. Unless specified otherwise, we set the hyperparameters to  $\beta_{\text{CM}} = 1$  and  $\beta_{\text{DM}} = 0.005$ . For embeddings, we use Gaussian initialization  $\mathcal{N}(0, D^{-1/2})$ . The EMA rate is set to  $\eta = 0.99$  and the embedding dimensionality is set to  $D = 256$ .

## B LATENT VARIABLE CLASSIFIER-FREE GUIDANCE

It is important to generate images corresponding to a given condition. In VQ-LCMD, the condition is incorporated directly into the prediction network through Adaptive Layer Normalization (Ba et al., 2016). The assumption here is that the network uses both the corrupted input and the condition to reconstruct the input. However, we often observe that VQ-LCMD generates outputs that are not correlated well with the condition. The reason is that the corrupted input contains rich information; therefore, the network can ignore the condition during training.

To improve the sample quality of conditional diffusion models, we employ the classifier-free guidance (Ho & Salimans, 2021). Essentially, it guides the sampling trajectories toward higher-density data regions. During training, we randomly drop 10% of the conditions and set the dropped conditions to the null token. During sampling, VQ-LCMD predicts the categorical variable  $\mathbf{x}$  as follows

$$\log P_{\theta}(\mathbf{x}|\mathbf{z}_t, \mathbf{y}; t) = (1 + \omega) \log P_{\theta}(\mathbf{x}|\mathbf{z}_t, \mathbf{y}; t) - \omega \log P_{\theta}(\mathbf{x}|\mathbf{z}_t; t), \quad (4)$$

where  $\omega \geq 0$  denotes the guidance scale and  $\mathbf{y}$  denotes the condition. Note that both terms on the right-hand side of Eq. (4) are parameterized by the same model. Figure 4 shows the effects of increasing the classifier-free guidance scale  $\omega$ .

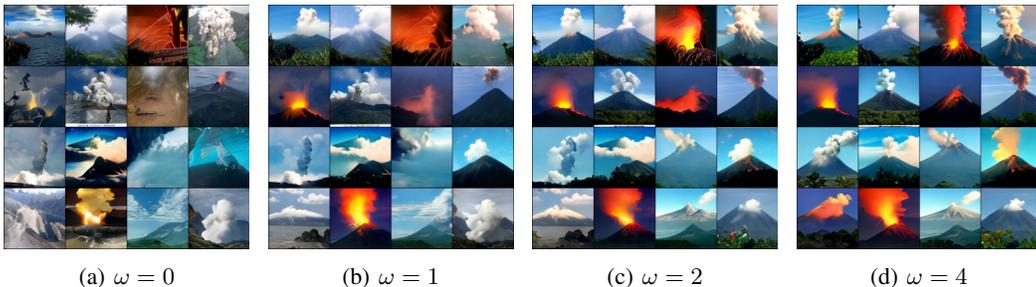


Figure 4: Generated samples of VQ-LCMD with  $\omega$  ranging from 0 to 4 on ImageNet.

### C ABLATION STUDIES

In this section, we provide additional ablation studies to further validate our motivations of VQ-LCMD.

#### C.1 PRETRAINED VS. LEARNABLE EMBEDDINGS

We evaluate the embedding vectors obtained by VQ-LCMD against those provided by the pretrained VQGAN on the LSUN Churches dataset. Figure 5 presents the magnitudes of these vectors and the distance matrices between embeddings. Interestingly, our method learns a structure that is quite similar to the pretrained embeddings. Learnable embeddings tend to have larger magnitudes.

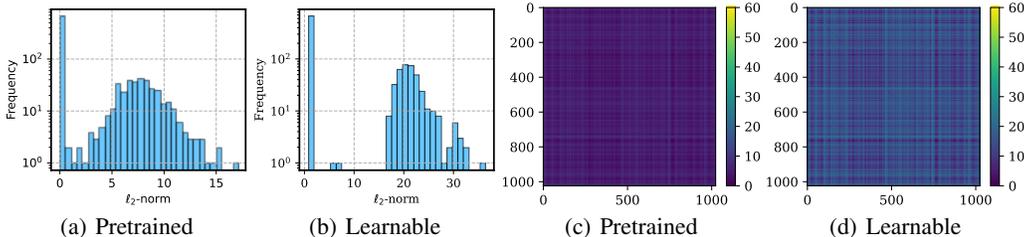


Figure 5: Visual representation of pretrained and learnable embedding vectors for the LSUN Churches dataset: (a) vector magnitudes for pretrained embeddings, (b) vector magnitudes for learnable embeddings, (c) distance matrix for pretrained embeddings and (d) distance matrix for learnable embeddings. For distance matrices, we compute the Euclidean distances between different embedding vectors.

#### C.2 NUMBER OF SAMPLING STEPS

We analyze the number of steps necessary to obtain high-fidelity samples. Table 4 presents the FID scores corresponding to different numbers of sampling steps. As expected, we observe a decrease in FID as the number of sampling steps increases. However, the improvement becomes marginal after reaching 50 steps. VQ-LCMD can accelerate the conventional diffusion models by a large margin, which is a notable advantage compared to ARs. In addition, we leverage the DDIM sampler to further reduce the number of sampling steps. Table 5 presents a comparison of VQ-LCMD using DDIM against MaskGIT (Besnier & Chen, 2023) (PyTorch implementation). The results show that VQ-LCMD achieves a better FID score than MaskGIT when the number of steps is extremely small, highlighting the advantage of our method.

#### C.3 DROPPING STRATEGIES

We explore three different strategies to drop tokens during training. One strategy involves linearly increasing the dropping ratio concerning the timestep (**linear**). In this scheme, early timesteps involve a small portion of tokens being dropped, while in later timesteps a higher proportion of tokens

Table 4: FID results for different numbers of sampling steps

Steps	5	10	15	20	50	100	200
Churches	19.38	10.24	7.81	6.80	5.43	5.20	<b>4.99</b>
Bedrooms	14.55	6.05	4.42	4.00	<b>3.86</b>	4.01	4.16
FFHQ	28.80	15.55	11.44	9.57	7.56	7.34	<b>7.25</b>

Table 5: FID results on ImageNet for different number of sampling steps

Steps	2	3	4	5	6	7	8
MaskGIT	97.83	46.43	20.29	10.95	7.74	<b>6.79</b>	<b>6.80</b>
VQ-LCMD	<b>17.83</b>	<b>10.57</b>	<b>8.51</b>	<b>7.87</b>	<b>7.56</b>	7.45	7.44

are dropped. Another strategy is to randomly select a ratio and drop the tokens according to this ratio (**rand\_drop**). Finally, a fixed dropping ratio  $0 \leq r \leq 1$  can be employed (**rand(r)**). Table 6 summarizes the results. VQ-LCMD performs the best with an appropriately chosen fixed dropping ratio.

Table 6: Ablation results on different dropping strategies

	FID↓	Precision↑	Recall↑
linear	9.12	<b>0.72</b>	0.41
rand_drop	8.44	0.71	0.42
rand (0.1)	7.81	0.71	0.43
rand (0.2)	<b>7.25</b>	<b>0.72</b>	<b>0.46</b>
rand (0.3)	8.45	0.70	0.43
rand (0.4)	9.89	0.70	0.41
rand (0.5)	9.11	<b>0.72</b>	0.41

#### C.4 WEIGHTING TERMS

We hypothesize that balancing the reconstruction loss and the diffusion loss is crucial to preventing embedding collapse. In VQ-LCMD, this is achieved by tuning the hyperparameter  $\beta_{DM}$ . Table 7 presents the FID results on FFHQ for various combinations of  $\beta_{CM}$  and  $\beta_{DM}$ . Adjusting these parameters alters the contributions of the diffusion loss and the consistency-matching loss in the objective function. As indicated in the table, when  $\beta_{DM}$  is relatively large, the model still suffers from embedding collapse.

#### C.5 EMBEDDING DIMENSIONALITY

Table 8 shows the influence of embedding dimensionality. We report the FID results on FFHQ when varying the embedding dimensionality. VQ-LCMD demonstrates consistent performance across various dimensionalities. As the dimensionality increases, the performance slightly decreases. VQ-LCMD achieves the best result when  $D = 128$ .

Table 7: Results on  $\beta_{CM}$  and  $\beta_{DM}$ 

$\beta_{CM}$	$\beta_{DM}$	FID ↓
0.01	0.01	175.46
0.01	1	173.28
1	1	54.10
1	0.01	8.26
1	0.005	7.25

Table 8: Embedding dimensionality

$D$	FID ↓
64	7.90
128	7.20
256	7.25
768	7.42
1024	7.38

### C.6 ABLATION STUDIES ON IMAGENET

We conduct ablation studies on ImageNet to examine the effects of classifier-free guidance weights and the number of sampling steps. Figure 6(a) shows the FID and IS metrics across various classifier-free guidance weight values. Additionally, Figure 6(b) presents the FID and IS results as we vary the number of sampling steps. There is a clear trade-off between fidelity represented by FID and quality represented by IS. VQ-LCMD achieves the best FID results when  $\omega = 1$ .

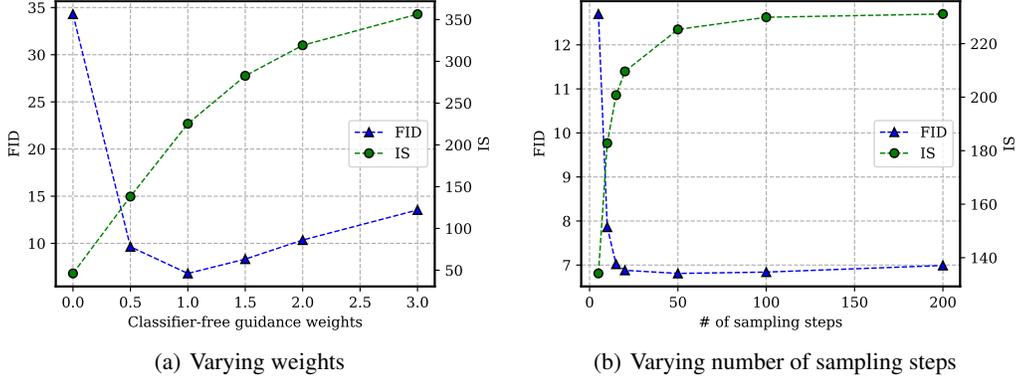


Figure 6: Ablation studies for FID vs IS on ImageNet when (a) varying classifier-free guidance weights and (b) varying number of sampling steps.

### D PSEUDOCODE

Algorithms 1 and 2 outline the training and sampling procedures of VQ-LCMD. For sampling, we discretize time  $t \in [0, 1]$  into  $N + 1$  points  $\{t_n\}_{n=0}^N$  such that they satisfy  $t_n < t_{n+1}$ ,  $t_0 = 0$ , and  $t_N = 1$ . Starting with Gaussian noise sampled from  $\mathbf{z}_{t_N} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , we sample  $\mathbf{z}_0$  through the ancestral sampling given by  $p_{\phi, \theta}(\mathbf{z}_{t_{n-1}} | \mathbf{z}_{t_n}; t_{n-1}, t_n)$ , which is defined in Eq. (1). Finally, the discrete output  $\mathbf{x}$  is obtained from the model  $P_{\theta}(\mathbf{x} | \mathbf{z}_0; 0)$ . Note that, unlike CSDMs, our model directly outputs the token probabilities for continuous input  $\mathbf{z}_{t_n}$  at timestep  $t_n$ .

### E ADDITIONAL SAMPLES

In this section, we present additional samples generated by VQ-LCMD. For unconditional image generation, Figures 7, 8, and 9 show the generated samples from VQ-LCMD trained on FFHQ, LSUN Churches, and LSUN Bedrooms, respectively. Figure 10 visualizes the conditional samples from ImageNet. All images are at a resolution of  $256 \times 256$ .

**Algorithm 1** Training

---

```

1: repeat
2:   Sample batch of  $\mathbf{x} \sim P(\mathbf{x})$ 
3:    $t \sim \mathcal{U}(0, 1)$ ;  $s \sim \mathcal{U}(0, t)$ ;  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ;
      $\mathbf{m}_{\text{RD}} \sim \{0, 1\}^M$ 
4:    $\mathbf{z}_t := \alpha_t \Psi_\phi(\mathbf{x}) \odot \mathbf{m}_{\text{RD}} + \sigma_t \epsilon$ 
5:    $\bar{\mathbf{z}}_s := \alpha_s \Psi_{\bar{\phi}}(\mathbf{x}) + (\sigma_s / \sigma_t)(\mathbf{z}_t - \alpha_t \Psi_\phi(\mathbf{x}))$ 
6:   Take gradient descent step on
7:      $\nabla_{\phi, \theta} \mathcal{L}(\mathbf{x}; \phi, \theta)$ 
8: until converged

```

---

**Algorithm 2** Sampling

---

```

1: Prepare
    $t_0 := 0 < t_1 < \dots < t_N := 1$  and
    $\mathbf{z}_{t_N} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $n = N, N - 1, \dots, 1$  do
3:    $\mathbf{z}_{t_{n-1}} \sim p_{\phi, \theta}(\mathbf{z}_{t_{n-1}} | \mathbf{z}_{t_n}; t_{n-1}, t_n)$ 
4: end for
5:  $\mathbf{x} \sim P_\theta(\mathbf{x} | \mathbf{z}_0; 0)$ 
6: return  $\mathbf{x}$ 

```

---



Figure 7: VQ-LCMD samples of unconditional image generation on FFHQ.



Figure 8: VQ-LCMD samples of unconditional image generation on LSUN Churches.

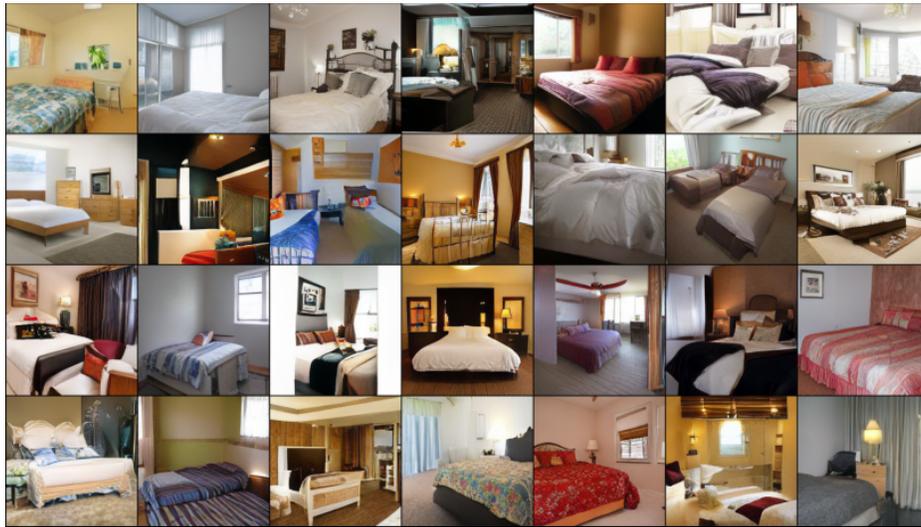


Figure 9: VQ-LCMD samples of unconditional image generation on LSUN Bedrooms.

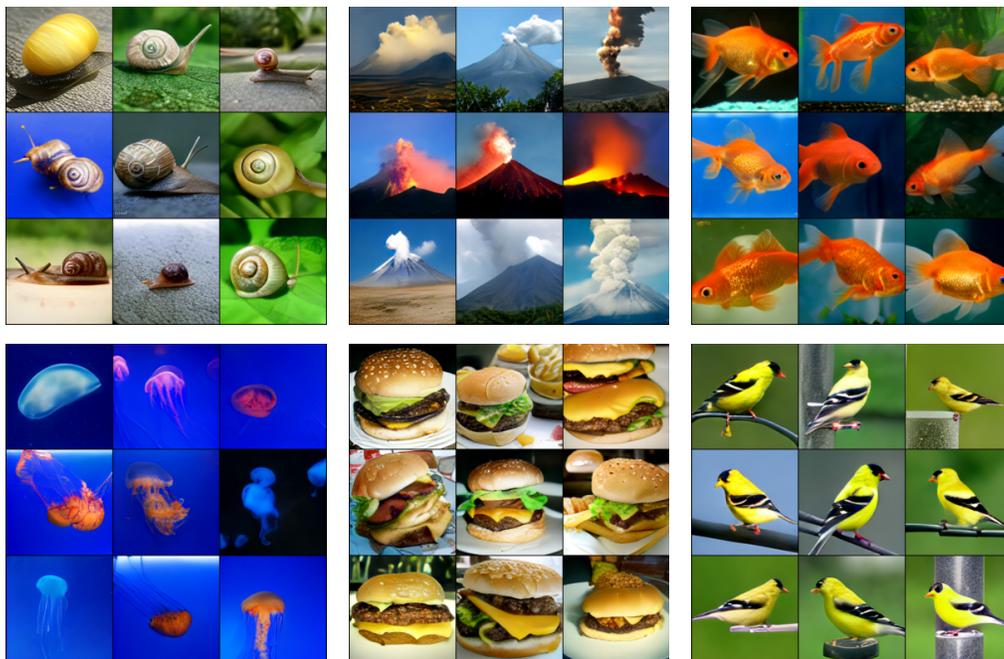


Figure 10: VQ-LCMD samples of conditional image generation on ImageNet  $256 \times 256$  for selected classes, including “snail”, “volcano”, “goldfish”, “jellyfish”, “cheeseburger”, “goldfinch”