

---

# Q-DM: An Efficient Low-bit Quantized Diffusion Model

---

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

1 Denoising diffusion generative models are capable of generating high-quality data,  
2 but suffers from the computation-costly generation process, due to a iterative noise  
3 estimation using full-precision networks. As an intuitive solution, quantization  
4 can significantly reduce the computational and memory consumption by low-  
5 bit parameters and operations. However, low-bit noise estimation networks in  
6 diffusion models (DMs) remain unexplored yet and perform much worse than the  
7 full-precision counterparts as observed in our experimental studies. In this paper,  
8 we first identify that the bottlenecks of low-bit quantized DMs come from a large  
9 distribution oscillation on activations and accumulated quantization error caused  
10 by the multi-step denoising process. To address these issues, we first develop a  
11 Timestep-aware Quantization (TaQ) method and a Noise-estimating Mimicking  
12 (NeM) scheme for low-bit quantized DMs (Q-DM) to effectively eliminate such  
13 oscillation and accumulated error respectively, leading to well-performed low-bit  
14 DMs. In this way, we propose an efficient Q-DM to calculate low-bit DMs by  
15 considering both training and inference process in the same framework. We evaluate  
16 our methods on popular DDPM and DDIM models. Extensive experimental results  
17 show that our method achieves a much better performance than the prior arts. For  
18 example, the 4-bit Q-DM theoretically accelerates the 1000-step DDPM by 7.8×  
19 and achieves a FID score of 5.17, on the unconditional CIFAR-10 dataset.

这两件事是  
一回事呢？

how?

20 

## 1 Introduction

21 Denoising diffusion models, also known as score-based generative models [10, 33, 35], have recently  
22 shown remarkable success in various generative tasks such as images [10, 35, 22], audio [21],  
23 video [31], and graphs [23]. These models have also demonstrated flexibility in downstream tasks,  
24 making them attractive for tasks such as super-resolution [26, 7] and image-to-image translation [29].  
25 Compared to Generative Adversarial Networks (GANs) [8], historically considered state-of-the-  
26 art, diffusion models have proven to be superior in terms of quality and diversity in most of these  
27 tasks and applications. The process of diffusion models involves gradually transforming real data  
28 into Gaussian noise, which is then reversed via a denoising process to generate real data [10, 40].  
29 However, such denoising process is time-consuming and involves iterating a neural network for  
30 noise estimation over thousands of timesteps, despite producing a significant amount of images.  
31 Therefore, researchers are actively working on accelerating this generation process to reduce its long  
32 iterative process and high inference cost for sample generation. To achieve this, one pipeline is to  
33 focus on sample trajectory learning, to develop faster sampling strategies [28, 22, 1]. While the  
34 other pipeline directly compresses and accelerates the noise estimation networks based on network  
35 quantization technology [30], which is particularly suitable for AI chips because of the low-bit  
36 parameters and operations. Prior post-training quantization (PTQ) methods [30, 19, 17] on diffusion  
37 models (DMs) or other neural networks directly compute quantized parameters based on pre-trained

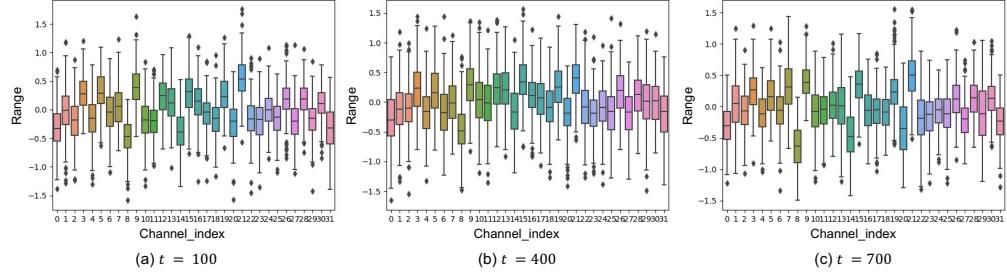


Figure 1: Studies on the activation distribution w.r.t. time-step. Per (output) channel activation ranges of the first attention block in diffusion model on different timestep. The boxplot visualizes key statistical measures for each channel, including the minimum and maximum values, the 2nd and 3rd quartiles, and the median.

38 full-precision models, which constrains the model performance to a sub-optimized level without  
39 fine-tuning. Furthermore, quantizing DMs based on PTQ methods to ultra-low bits (e.g., 4 bits or  
40 lower) is ineffective and suffers from a significant performance reduction.

41 Differently, quantization-aware training (QAT) [16, 18] methods perform quantization during back  
42 propagation and generally achieve a less performance drop with a higher compression rate than PTQ.  
43 For instance, QAT has been shown to be effective for CNNs [5, 18] and ViTs [16, 18] and BERT [24].  
44 However, QAT methods for low-bit quantization of diffusion models remain largely unexplored.  
45 Therefore, we first build a low-bit quantized DM baseline, a straightforward yet effective solution  
46 based on common techniques [5]. Our experimental studies reveal that the severe performance drop  
47 of low-bit quantized DMs, such as PTQ [30] and baseline [5], lies in the activation distribution  
48 oscillation and quantization error accumulation caused by the denoising process.

49 As shown in Fig. 1, the output distribution of the noise estimation network at each time step can  
50 differ significantly, resulting in activation distribution oscillation. Particularly, the distribution of  
51 activation in a specific layer varies significantly across different timesteps during training. We also  
52 observe that errors between full-precision activations and quantized activations gradually accumulate  
53 across timesteps during the sampling process (inference), making it harder to produce well-performed  
54 quantized DMs. 观察2：量化误差积累

观察1：分布振荡 ACT  
不同 time step 分布  
范围不同。

55 Drawing on the aforementioned insights, we propose a Timestep-aware Quantization (TaQ) method  
56 to address the oscillating distribution issue. By smoothing out these fluctuations and introducing  
57 more precise scaling factors into activations, we effectively enhance the performance of the low-bit  
58 quantized DMs. We further design a new training scheme for quantized DMs, dubbed Noise-  
59 estimating Mimicking (NEM), which can reduce the accumulated errors and promote the performance  
60 of quantized DMs based on the knowledge of full-precision counterparts. In this way, we achieve a  
61 new QAT method for low-bit quantized DM (Q-DM) via incorporating all the explorations (see the  
62 overview in Fig. 2). Overall, the contributions of this paper can be summarized as follows:

- 63 • To the best of our knowledge, we proposed the first QAT method towards efficient low-bit  
64 DMs, dubbed Q-DM, by fully considering both training and inference process in the same  
65 framework. 在同一个框架内考虑训练与测试  
是什么意思？
- 66 • We introduce a Timestep-aware Quantization (TaQ) method to mitigate activation distribution  
67 oscillation caused by the random-sampled timestep in the training process. We develop a  
68 Noise-estimating Mimicking (NEM) scheme to reduce accumulated errors, by which the  
69 Q-DMs are able to achieve comparable performance as the full-precision counterparts.
- 70 • Extensive experiments on the CIFAR-10 and ImageNet datasets show that our Q-DM  
71 outperforms the baseline and 8-bit PTQ method by a large margin, and achieves comparable  
72 performances as the full-precision counterparts with a considerable acceleration rate.

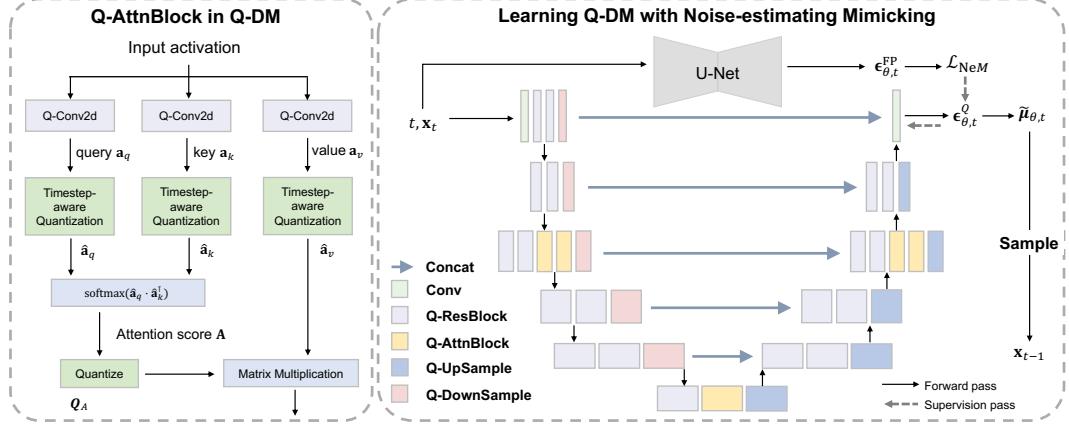


Figure 2: Overview of the proposed Q-DM framework. We introduce the timestep-aware quantization in an architecture perspective and a noise imitation training scheme incorporated in the optimization process. From left to right, we respectively show the detailed architecture of single Q-AttnBlock in Q-DM and the training framework of Q-DM.

这东西是所有 diffusion model  
都有的吗？

## 73 2 Related Work

74 **Network Quantization.** Quantizing neural networks (QNNs) often possess low-bit (1~4-bit) weights  
 75 and activations to accelerate the model inference and save the memory usage. Specifically, ternary  
 76 weights are introduced to reduce the quantization error in TWN [15]. DoReFa-Net [41] exploits  
 77 convolution kernels with low bit-width parameters and gradients to accelerate both the training and  
 78 inference. TTQ [42] uses two full-precision scaling coefficients to quantize the weights to ternary  
 79 values. Zhuang *et al.* [43] present a 2 ~ 4-bit quantization scheme using a two-stage approach  
 80 to alternately quantize the weights and activations, which provides an optimal trade-off among  
 81 memory, efficiency, and performance. Jung *et al.* [12] parameterize the quantization intervals and  
 82 obtain their optimal values by directly minimizing the task loss of the network and also the accuracy  
 83 degeneration with further bit-width reduction. ZeroQ [2] supports both uniform and mixed-precision  
 84 quantization by optimizing for a distilled dataset, which is engineered to match the statistics of batch  
 85 normalization across different layers of the network. Xie *et al.* [39] introduces transfer learning into  
 86 network quantization to obtain an accurate low-precision model by utilizing the Kullback-Leibler  
 87 (KL) divergence. PWLQ [6] enables accurate approximation for tensor values that have bell-shaped  
 88 distributions with long tails and finds the entire range by minimizing the quantization error.

89 **Diffusion Model.** The high cost of denoising through networks and the long iterative process  
 90 make it difficult to implement diffusion models widely. To accelerate diffusion probabilistic models  
 91 (DMs) [10], previous research has focused on finding shorter sampling trajectories while maintaining  
 92 DM performance. Wavegrad [3] introduces grid search, which finds an effective trajectory with only  
 93 six timesteps, but this approach cannot be generalized for longer trajectories due to its exponentially  
 94 growing time complexity. Watson *et al.* [38] model the trajectory searching as a dynamic programming  
 95 problem. Song *et al.* [34] construct non-Markovian diffusion processes that lead to the same training  
 96 objective, but whose reverse process can be much faster to sample from. For DMs with continuous  
 97 timesteps, Song *et al.* [33, 35] have formulated the DM in the form of an ordinary differential  
 98 equation (ODE) and improved sampling efficiency by using faster ODE solvers. Jolicoeur-Martineau  
 99 *et al.* [11] have introduced an advanced SDE solver to accelerate the reverse process via an adaptively  
 100 larger sampling rate. Analytic-dpm [1] has estimated variance and KL divergence using the Monte  
 101 Carlo method and a pretrained score-based model with derived analytic forms that are simplified  
 102 from the score-function. In addition to those training-free methods, Luhman & Luhman [20] have  
 103 compressed the reverse denoising process into a single-step model, while San-Roman *et al.* [28]  
 104 has dynamically adjusted the trajectory during inference. However, implementing these methods  
 105 requires additional training after obtaining a pretrained DM, which makes them less desirable in most  
 106 situations. In summary, all these DM acceleration methods can be categorized as finding effective  
 107 sampling trajectories.

108 Unlike prior works, we demonstrate that diffusion models can be accelerated by compressing the  
 109 network in each noise estimating iteration, which is orthogonal with the fast sampling methods  
 110 mentioned above. To the best of our knowledge, this is the first study to explore low-bit quantized  
 111 diffusion models in a quantization-aware training (QAT) manner.

### 112 3 Background and Challenge

#### 113 3.1 Diffusion Models

114 **Forward process.** Let  $\mathbf{x}_0$  be a sample from the data distribution  $\mathbf{x}_0 \sim q(\mathbf{x})$ . A forward diffusion  
 115 process adds Gaussian noise to the sample for  $T$  times, resulting in a sequence of noisy samples  
 116  $\mathbf{x}_1, \dots, \mathbf{x}_T$  as:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad (1)$$

117 where  $\beta_t \in (0, 1)$  is the variance schedule and controls the strength of the Gaussian noise in each  
 118 step. The forward diffusion process satisfies the Markov property since each step relies solely on the  
 119 preceding step. Additionally, as the number of steps increases towards infinity ( $T \rightarrow \infty$ ), the final  
 120 state  $\mathbf{x}_T$  converges to an isotropic Gaussian distribution. A notable property of the forward process is  
 121 that it admits sampling  $\mathbf{x}_t$  at an arbitrary timestep  $t$  in closed form as:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}). \quad (2)$$

122 **Reverse process.** To generate a sample from a Gaussian noise input  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  using diffusion  
 123 models, the forward process is reversed. However, since the actual reverse conditional distribution  
 124  $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$  is unknown, diffusion models use a learned conditional distribution  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$  that ap-  
 125 proximates the real reverse conditional distribution with a Gaussian distribution. This approximation  
 126 is expressed as:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_{\theta,t}(\mathbf{x}_t), \tilde{\beta}_t \mathbf{I}). \quad (3)$$

127 By using the re-parameterization trick presented in [10], it becomes possible to derive the mean  
 128  $\tilde{\mu}_{\theta,t}(\mathbf{x}_t)$  and  $\tilde{\beta}_t \mathbf{I}$  as follows:

$$\begin{aligned} \tilde{\mu}_{\theta,t}(\mathbf{x}_t) &= \frac{1}{\sqrt{\alpha_t}} (\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta,t}), \\ \tilde{\beta}_t &= \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t, \end{aligned} \quad (4)$$

129 where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$  and  $\epsilon_\theta$  is a function approximator intended to predict  $\epsilon$  from  
 130  $\mathbf{x}_t$  [10].

131 **Training.** At training time, the goal of optimization is to minimize the negative log-likelihood, *i.e.*,  
 132  $-\log p_\theta(\mathbf{x}_0)$ . With variational inference, a lower bound of it could be found, denoted as  $L_{\text{VLB}}$ :

$$L_{\text{VLB}} = \mathbb{E}_{q(\mathbf{x}_{0:T})} [\log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})}] \geq -\log p_\theta(\mathbf{x}_0). \quad (5)$$

133 It is found in [10] that using a simplified loss function to  $L_{\text{VLB}}$  often obtains better performance:

$$L_{\text{simple}} = \mathbb{E}_{t, \mathbf{x}_0, \epsilon_t} [\|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2]. \quad (6)$$

134 **Sampling.** At inference time, a Gaussian noise tensor  $\mathbf{x}_T$  is sampled and is denoised by repeatedly  
 135 sampling the reverse distribution  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$ .  $\tilde{\mu}_{\theta,1}(\mathbf{x}_1)$  is taken as the final generation result, with  
 136 no noise added in the final denoising step.

#### 137 3.2 Quantization

138 Given an  $N$ -layer CNN model, we denote its weight set as  $\mathbf{W} = \{\mathbf{w}^n\}_{n=1}^N$  and input feature map set  
 139 as  $\mathbf{A} = \{\mathbf{a}_{\text{in}}^n\}_{n=1}^N$ . The  $\mathbf{w}^n \in \mathbb{R}^{C_{\text{out}}^n \times C_{\text{in}}^n \times K^n \times K^n}$  and  $\mathbf{a}_{\text{in}}^n \in \mathbb{R}^{C_{\text{in}}^n \times W_{\text{in}}^n \times H_{\text{in}}^n}$  are the convolutional  
 140 weight and the input feature map in the  $n$ -th layer, where  $C_{\text{in}}^n$ ,  $C_{\text{out}}^n$  and  $K^n$  respectively stand for  
 141 input channel number, output channel number and the kernel size. Also,  $W_{\text{in}}^n$  and  $H_{\text{in}}^n$  are the width

① why DDPM?  
② 同一模型？

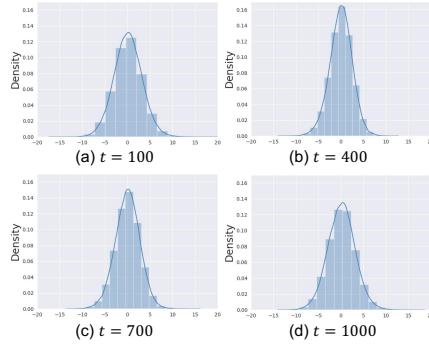


Figure 3: Input activation distribution of the first Q-AttnBlock in diffusion model on different timestep with a model trained on CIFAR-10 [13] by DDPM.

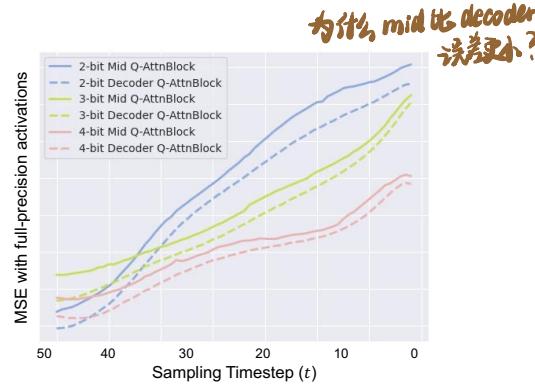


Figure 4: Distance between the outputs of the full-precision model and different bit-width baseline models trained on CIFAR-10 [13] by DDIM with 100 sampling steps.

142 and height of the feature maps. Then, the convolutional outputs  $\mathbf{a}_{\text{out}}^n$  can be technically formulated  
143 as:

$$\mathbf{a}_{\text{out}}^n = \mathbf{w}^n \otimes \mathbf{a}_{\text{in}}^n, \quad (7)$$

144 where  $\otimes$  represents the convolution operation. Herein, we omit the non-linear function for simplicity.  
145 Quantized neural network intends to represent  $\mathbf{w}^n$  and  $\mathbf{a}^n$  in a low-bit format such that the float-point  
146 convolutional outputs can be approximated as:

$$\begin{aligned} \hat{\mathbf{w}}^n &= s^{w^n} \circ Q(\mathbf{w}^n) = s^{w^n} \cdot [\text{clip}(\mathbf{w}^n / s^{w^n}, -2^{b-1}, 2^{b-1} - 1)] \\ \hat{\mathbf{a}}_{\text{in}}^n &= s^{a_{\text{in}}^n} \cdot Q(\mathbf{a}_{\text{in}}^n) = s^{a_{\text{in}}^n} \cdot [\text{clip}(\mathbf{a}_{\text{in}}^n / s^{a_{\text{in}}^n}, -2^{b-1}, 2^{b-1} - 1)] \\ \hat{\mathbf{a}}_{\text{out}}^n &= \hat{\mathbf{a}}_{\text{in}}^n \otimes \hat{\mathbf{w}}^n \approx s^{a_{\text{in}}^n} \cdot s^{w^n} \circ [Q(\mathbf{w}^n) \odot Q(\mathbf{a}_{\text{in}}^n)], \end{aligned} \quad (8)$$

147 where  $\circ$  denotes the channel-wise multiplication,  $\odot$  denotes the efficient GEMM operations, and  
148  $s^{w^n} = \{s_1^{w^n}, s_2^{w^n}, \dots, s_{C_{\text{out}}}^{w^n}\} \in \mathbb{R}_+^{C_{\text{out}}}$  is known as the channel-wise scaling factor vector [25] to  
149 mitigate the output gap between Eq. (7) and its approximation of Eq. (8). Meanwhile, we use the  
150 layer-wise quantization for input activations and the scaling factor of activations  $s^{a_{\text{in}}^n} \in \mathbb{R}_+$  is a  
151 scalar.

### 152 3.3 Challenge Analysis

有两类训练步骤。

153 Here we identify two major challenges on low-bit DMs, specific to the multi-step inference process  
154 and random-sampled-step training process of diffusion models. Namely, we investigate on the  
155 distribution oscillation of the activations, and the accumulated quantization error resulted from the  
156 multi-step denoising process.

good insight

157 **Activation distribution oscillation.** To understand the distribution change of diffusion models, we  
158 investigate the activation distribution, w.r.t. timestep in the training process. Theoretically, if the  
159 distribution changes w.r.t. timestep, it would be difficult to implement previous QAT methods. We  
160 analyze the overall activation distributions of the noise estimation network, as shown in Fig. 3. We can  
161 observe that at different timesteps, the corresponding activation distributions have large discrepancies,  
162 e.g., Fig. 3(a) v.s. Fig. 3(b), which makes previous QAT methods [16] in-applicable for multi-timestep  
163 models, i.e., diffusion models.

一个基本问题：不同  
time step 对应的模  
型是同一个吗？

为什么？只有同一个模型处理不同 time step 才会有这样的观察点？

是的，但是 practice 3?

164 **Quantization error accumulation.** Quantization of a noise estimation network introduces distur-  
165 bances to the weights and activations, resulting in errors in each layer's output. Previous studies [4]  
166 have found that these errors tend to accumulate across layers, making it more challenging to quantize  
167 deeper neural networks. In the case of diffusion models (DMs), at each time step  $t$ , the input of the  
168 model ( $\mathbf{x}_{t-1}$ ) is obtained from the model's output at the previous time step  $t$ , i.e.,  $\mathbf{x}_t$ . As depicted  
169 in Fig. 4, the MSE distance, representing the quantization error of low-bit quantized DMs, exhibits  
170 a noticeable growth along with the decrease of sampling timestep. This implies that as the denois-  
171 ing process moves towards later timestep, the accumulation of quantization errors becomes more  
172 prominent.

173 4 The Proposed Q-DM

②目前常用的 V-Net 有 attention block 吗?

174 4.1 Timestep-aware Quantization

175 To tackle the distribution oscillation in the training process, we first introduce the quantized attention  
176 block that efficiently takes into account the timestep. This structure allows for the numerical  
177 analysis of activation ranges across different timesteps and mitigates distribution oscillation of low-bit  
178 quantized DMs. We recall the quantization in the attention block based on Eq. (8), which is formulated  
179 as:

$$\begin{aligned}\hat{\mathbf{a}}_q(\mathbf{x}_t, t) &= s^{a_q(\mathbf{x}_t, t)} \cdot Q(\mathbf{a}_q(\mathbf{x}_t, t)), \quad \hat{\mathbf{a}}_k(\mathbf{x}_t, t) = s^{a_k(\mathbf{x}_t, t)} \cdot Q(\mathbf{a}_k(\mathbf{x}_t, t)) \\ \mathbf{A}(\mathbf{x}_t, t) &= \text{softmax}[(\hat{\mathbf{a}}_q(\mathbf{x}_t, t) \cdot \hat{\mathbf{a}}_k(\mathbf{x}_t, t)^\top) / \sqrt{d}], \\ \hat{\mathbf{A}}(\mathbf{x}_t, t) &= s^{A(\mathbf{x}_t, t)} \cdot Q(\mathbf{A}(\mathbf{x}_t, t)), \\ \mathbf{a}_{\text{out}}(\mathbf{x}_t, t) &= \hat{\mathbf{A}}(\mathbf{x}_t, t) \cdot \hat{\mathbf{a}}_v(\mathbf{x}_t, t)^\top,\end{aligned}\tag{9}$$

180 where  $\mathbf{A}$  is the attention score.

181 In the  $i$ -th mini-batch, the timestep is represented as  $\{t_1, \dots, t_{b_i}\}$ , where  $b_i$  is the batch size of the  
182  $i$ -th batch. We denote  $i \in \{1, \dots, B\}$ , and  $B$  is the number of batches. Therefore, we calculate the  
183 timestep-aware distribution divergence for the query activation  $\mathbf{a}_q$  as:

这里本质上是做了个  
time-step aware 的统计元素  
中心化吗?

$$\begin{aligned}\gamma_{q;t} &= \sum_{i=1}^B \frac{1}{b_i} \sum_{j=1}^{b_i} \underbrace{\mathbf{a}_q(\mathbf{x}_{t_j}, t_j)}_{\text{??不就是指的这个吗?}}, \quad \text{④这到底是怎么确定的? 不会随着训练更新吗?} \\ \sigma_{q;t}^2 &= \sum_{i=1}^B \frac{1}{b_i} \sum_{j=1}^{b_i} [\mathbf{a}_q(\mathbf{x}_{t_j}, t_j) - \gamma_{q;t}]^2,\end{aligned}\tag{10}$$

184 where  $\gamma_{q;t}$  and  $\sigma_{q;t}^2$  are statistical mean and variance of query activation  $\mathbf{a}_q$ . And the calculation of  
185 the key activation  $\mathbf{a}_k$  is likewise.

186 Based on such statistical results, the query and key activations in each specific timestep are smoothed  
187 as:

$$\begin{aligned}\tilde{\mathbf{a}}_q(\mathbf{x}_t, t) &= [\mathbf{a}_q(\mathbf{x}_t, t) - \gamma_{q;t}] / \sqrt{\sigma_{q;t}^2 + \psi} \\ \tilde{\mathbf{a}}_k(\mathbf{x}_t, t) &= [\mathbf{a}_k(\mathbf{x}_t, t) - \gamma_{k;t}] / \sqrt{\sigma_{k;t}^2 + \psi},\end{aligned}\tag{11}$$

188 where  $\psi$  is constant to avoid 0 denominator. With the above timestep-aware smoothing process, we  
189 formulate our timestep-aware quantization as:

$$\begin{aligned}\hat{\mathbf{a}}_q(\mathbf{x}_t, t) &= s^{a_q(\mathbf{x}_t, t)} \cdot \text{TaQ}(\mathbf{a}_q(\mathbf{x}_t, t)), \quad \hat{\mathbf{a}}_k(\mathbf{x}_t, t) = s^{a_k(\mathbf{x}_t, t)} \cdot \text{TaQ}(\mathbf{a}_k(\mathbf{x}_t, t)) \\ \mathbf{A}(\mathbf{x}_t, t) &= \text{softmax}[(\hat{\mathbf{a}}_q(\mathbf{x}_t, t) \cdot \hat{\mathbf{a}}_k(\mathbf{x}_t, t)^\top) / \sqrt{d}], \\ \hat{\mathbf{A}}(\mathbf{x}_t, t) &= s^{A(\mathbf{x}_t, t)} \cdot \text{TaQ}(\mathbf{A}(\mathbf{x}_t, t)), \\ \mathbf{a}_{\text{out}}(\mathbf{x}_t, t) &= \hat{\mathbf{A}}(\mathbf{x}_t, t) \cdot \hat{\mathbf{a}}_v(\mathbf{x}_t, t)^\top,\end{aligned}\tag{12}$$

190 in which  $\text{TaQ}(\mathbf{x}) = [\text{clip}([\mathbf{x} - \gamma_{*,t}] / [s^* \cdot \sqrt{\sigma_{*,t}^2 + \psi}], -2^{b-1}, 2^{b-1} - 1)]$ . The smoothed activations  
191 are less sensitive to the random sampled timestep in the training process and the timestep-aware  
192 quantization, to some extent, dismisses the distribution oscillation phenomenon.

193 4.2 Noise-estimating Mimicking

194 To mitigate the negative impact of quantization error accumulation on the training of a quantized  
195 DM  $\theta^Q$ , a full-precision DM, denoted as  $\theta^{\text{FP}}$ , is incorporated into the training process to facil-  
196 itate the learning objective. Following [10], with  $p_{\theta^Q}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_{\theta^Q, t}(\mathbf{x}_t), \tilde{\beta}_t \mathbf{I})$  and  
197  $p_{\theta^{\text{FP}}}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_{\theta^{\text{FP}}, t}(\mathbf{x}_t), \tilde{\beta}_t \mathbf{I})$ , we can write:

$$L_{t-1} = \mathbb{E}_q \left[ \frac{1}{2\tilde{\beta}_t} \|\boldsymbol{\mu}_{\theta^{\text{FP}}}(\mathbf{x}_t, t) - \boldsymbol{\mu}_{\theta^Q}(\mathbf{x}_t, t)\|^2 \right] + C,\tag{13}$$

Table 1: Evaluating the components of Q-DM based on 50-step DDIM sampler with  $32 \times 32$  generating resolution on CIFAR-10 [13]. “#Bits” denotes bit-width of weights and activations

Method	#Bits	FID $\downarrow$	IS $\uparrow$	#Bits	FID $\downarrow$	IS $\uparrow$	#Bits	FID $\downarrow$	IS $\uparrow$
Full-precision	32-32	4.67	9.27	-	-	-	-	-	-
PTQ4DM	8-8	18.02	8.87	-	-	-	-	-	-
Baseline (LSQ [5])	4-4	10.22	8.91	3-3	13.24	8.88	2-2	18.74	8.65
+TaQ	4-4	9.25	8.95	3-3	11.19	8.91	2-2	16.83	8.71
+NeM	4-4	8.98	8.92	3-3	11.02	8.90	2-2	16.97	8.79
<b>+TaQ+NeM (Q-DM)</b>	<b>4-4</b>	<b>6.89</b>	<b>8.96</b>	<b>3-3</b>	<b>9.07</b>	<b>8.98</b>	<b>2-2</b>	<b>15.26</b>	<b>8.86</b>

198 where  $C$  is a constant that does not depend on  $\theta^Q$  or  $\theta^{FP}$ . As in Eq. (13), we aim to compel the  
199 quantized model to replicate the noise estimation capability of the full-precision model. Further, by  
200 re-parameterizing Eq. (2) as  $\mathbf{x}_t(\mathbf{x}_0, \epsilon) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$  for  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and following the  
201 formulation in [10], which utilizes the formula for the posterior of the forward process, we can derive  
202 that:

$$L_{t-1} - C = \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[ \frac{1}{2\tilde{\beta}_t} \|\mu_{\theta^{FP}}(\mathbf{x}_t(\mathbf{x}_0, \epsilon), t) - \mu_{\theta^Q}(\mathbf{x}_t(\mathbf{x}_0, \epsilon), t)\|^2 \right], \quad (14)$$

203 where  $\mu_{\theta^{FP}}(\mathbf{x}_t(\mathbf{x}_0, \epsilon), t)$  and  $\mu_{\theta^Q}(\mathbf{x}_t(\mathbf{x}_0, \epsilon), t)$  are parameterized as:

$$\begin{aligned} \mu_{\theta^{FP}}(\mathbf{x}_t, t) &= \tilde{\mu}_t(\mathbf{x}_t, \frac{1}{\sqrt{\bar{\alpha}_t}} [\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta^{FP}}(\mathbf{x}_t)]) = \frac{1}{\alpha_t} (\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta^{FP}}(\mathbf{x}_t, t)), \\ \mu_{\theta^Q}(\mathbf{x}_t, t) &= \tilde{\mu}_t(\mathbf{x}_t, \frac{1}{\sqrt{\bar{\alpha}_t}} [\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_{\theta^Q}(\mathbf{x}_t)]) = \frac{1}{\alpha_t} (\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta^Q}(\mathbf{x}_t, t)), \end{aligned} \quad (15)$$

204 where  $\theta^Q$  and  $\theta^{FP}$  are the noise estimated by the quantized DM and full-precision counterpart. In  
205 Eq. (15),  $\epsilon_\theta$  is a function approximator intended to predict  $\epsilon$  from  $\mathbf{x}_t$ . Therefore, Eq. (14) is simplified  
206 to:

$$L_{t-1} = \mathbb{E}_{\mathbf{x}_0, \epsilon} \left[ \frac{\beta_t^2}{2\tilde{\beta}_t \alpha_t (1 - \bar{\alpha}_t)} \|\epsilon_{\theta^{FP}}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) - \epsilon_{\theta^Q}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2 \right]. \quad (16)$$

207 With the aforementioned derivation and parameterization, we have the final objective of our noise-  
208 estimating imitation, which is formulated as:

$$\begin{aligned} \arg \min_{\theta^Q} L_{\text{NeM}}(\theta^Q, \theta^{FP}) &\quad \text{⑤正常的训练目标是怎样的？我们可以理解成这个类似的吗} \\ &:= \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\|\epsilon_{\theta^{FP}}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) - \epsilon_{\theta^Q}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2], \quad \text{⑥DDPM和DDIM不是训练流程吗？也能拍代码吗？} \end{aligned} \quad (17)$$

## 209 5 Experiments

210 In this section, we evaluate the proposed Q-DM framework on several popular diffusion models (*i.e.*,  
211 **DDPM** [10] and **DDIM** [32]) for unconditional image generation. To the best of our knowledge,  
212 there is no published work done on low-bit quantized diffusion models at this point, so we report  
213 LSQ [5] as a baseline. Experiments show our approach can achieve competitive generation quality to  
214 the full-precision scenario on all experimental settings under low-bit quantization.

### 215 5.1 Datasets and Implementation Details

216 We evaluate our method on two datasets including **32** $\times$ **32** generating size in CIFAR-10 [13] and  
217 **64** $\times$ **64** generating size in ImageNet [14]. For the CIFAR-10 [13] and ImageNet [14] datasets, we  
218 use the DDIM [32] sampler with 50/100 sampling timesteps and DDPM [10] with 1000 sampling  
219 timesteps. All the training settings are the same as DDPM [10]. For DDIM sampler, we set  $\eta$  in  
220 DDIM [32] as 0.5 for the best performance. We evaluate the performance of our method using  
221 FID [9] and Inception Score (IS) [27] on both CIFAR-10 [13] and ImageNet [14] datasets. We set the  
222 training timestep  $T = 1000$  for all experiments, following [10]. We set the forward process variances

Table 2: Experiment on 2/3/4-bit quantized diffusion models generating CIFAR-10 [13] image or ImageNet [14] image. “#Bits” denotes the bit-width of weights/activations. “Reso.” represents the generating resolution.

Model	Dataset & Reso.	Step	Method	#Bits	Size <sub>(MB)</sub>	OPs <sub>(G)</sub>	FID ↓	IS ↑
DDIM	CIFAR-10 32×32	50	Full-precision	32/32	4.47	390.4	4.67	9.27
			PTQ4DM [6]	8/8	1.12	99.5	18.02	8.87
			Baseline	4/4	0.56	49.9	10.22	8.91
			<b>Q-DM</b>	4/4	0.56	49.9	<b>6.89</b>	<b>8.96</b>
			Baseline	3/3	0.28	25.1	13.24	8.88
		100	<b>Q-DM</b>	3/3	0.28	25.1	<b>9.07</b>	<b>8.98</b>
			Baseline	2/2	0.14	12.6	18.74	8.65
			<b>Q-DM</b>	2/2	0.14	12.6	<b>15.26</b>	<b>8.86</b>
			Full-precision	32/32	4.47	780.7	4.16	9.32
			PTQ4DM [6]	8/8	1.12	199.0	14.18	9.31
DDPM	CIFAR-10 32×32	1000	Baseline	4/4	0.56	99.8	9.02	8.95
			<b>Q-DM</b>	4/4	0.56	99.8	<b>5.12</b>	<b>9.21</b>
			Baseline	3/3	0.28	50.1	12.24	8.90
			<b>Q-DM</b>	3/3	0.28	50.1	<b>8.12</b>	<b>8.94</b>
			Baseline	2/2	0.14	25.2	16.99	8.74
		50	<b>Q-DM</b>	2/2	0.14	25.2	<b>14.31</b>	<b>8.77</b>
			Full-precision	32/32	4.47	7807.2	3.17	9.46
			PTQ4DM [6]	8/8	1.12	1990.0	7.10	9.55
			Baseline	4/4	0.56	997.7	9.11	8.96
			<b>Q-DM</b>	4/4	0.56	997.7	<b>5.17</b>	<b>9.15</b>
DDPM	ImageNet 64×64	100	Baseline	3/3	0.28	501.0	12.28	8.91
			<b>Q-DM</b>	3/3	0.28	501.0	<b>8.14</b>	<b>8.93</b>
			Baseline	2/2	0.14	252.0	16.93	8.72
			<b>Q-DM</b>	2/2	0.14	252.0	<b>14.35</b>	<b>8.76</b>
		50	Full-precision	32/32	4.47	390.4	20.57	15.72
			PTQ4DM [6]	8/8	1.12	99.5	25.87	14.99
			Baseline	4/4	0.56	49.9	24.78	15.37
			<b>Q-DM</b>	4/4	0.56	49.9	<b>20.02</b>	<b>15.68</b>
			Baseline	3/3	0.28	25.1	26.35	15.24
DDPM	ImageNet 64×64	1000	<b>Q-DM</b>	3/3	0.28	25.1	<b>22.19</b>	<b>15.32</b>
			Baseline	2/2	0.14	12.6	32.43	14.66
			<b>Q-DM</b>	2/2	0.14	12.6	<b>28.42</b>	<b>15.03</b>
			Full-precision	32/32	4.47	780.7	19.70	15.98
			PTQ4DM [6]	8/8	1.12	199.0	24.92	15.52
		50	Baseline	4/4	0.56	99.8	24.46	15.51
			<b>Q-DM</b>	4/4	0.56	99.8	<b>19.56</b>	<b>15.92</b>
			Baseline	3/3	0.28	50.1	26.23	15.42
			<b>Q-DM</b>	3/3	0.28	50.1	<b>21.97</b>	<b>15.92</b>
			Baseline	2/2	0.14	25.2	31.19	14.89
DDPM	ImageNet 64×64	1000	<b>Q-DM</b>	2/2	0.14	25.2	<b>27.94</b>	<b>14.99</b>
			Full-precision	32/32	4.47	7807.2	18.98	16.63
			PTQ4DM [6]	8/8	1.12	1990.0	22.32	15.31
			Baseline	4/4	0.56	997.7	22.91	15.29
			<b>Q-DM</b>	4/4	0.56	997.7	<b>18.52</b>	<b>16.72</b>
		50	Baseline	3/3	0.28	501.0	24.75	15.11
			<b>Q-DM</b>	3/3	0.28	501.0	<b>20.21</b>	<b>16.17</b>
			Baseline	2/2	0.14	252.0	29.33	14.87
			<b>Q-DM</b>	2/2	0.14	252.0	<b>25.62</b>	<b>15.48</b>

223 to constants increasing linearly from  $\beta_1 = 1e - 4$  to  $\beta_T = 0.02$ . To represent the reverse process, we  
224 use a U-Net backbone, following [10, 32]. Parameters are shared across time, which is specified to  
225 the network using the Transformer sinusoidal position embedding [36]. We use self-attention at the  
226  $16 \times 16$  feature map resolution [36, 37].

## 227 5.2 Ablation Study

228 We give quantitative results of the proposed TaQ and NeM in Tab. 1. As can be seen, the low-bit  
229 quantized DM baseline [5] suffers a severe performance drop on image generation task compared  
230 with full-precision DMs (5.55, 8.57, and 14.07 performance gap in terms of FID score with 4/3/2-bit,  
231 respectively). TaQ and NeM improve the performance of generation when used alone. For example,  
232 the 4-bit quantized DM baseline with TaQ and NeM introduced separately achieves 0.97 and 1.24 FID  
233 score decrease, respectively.

234 Moreover, the two techniques further boost the performance considerably when combined together.  
235 For instance, when combining the TaQ and NeM together, the performance of 4/3/2-bit quantized  
236 DMs improvement achieves 3.33, 4.07, and 3.48 respectively. To conclude, the two techniques can  
237 promote each other to improve Q-DM and close the performance gap between low-bit quantized  
238 DMs and full-precision counterpart.

## 239 5.3 Main Results

240 The experimental results are shown in Tab. 2. We compare our method with 4/3/2-bit baseline [5]  
241 based on the same frameworks for the task of unconditional image generation with the CIFAR-10 [13]  
242 and ImageNet [14] dataset. We also report the classification performance of the 8-bit PTQ method,  
243 i.e., PTQ4DM [30]. We firstly evaluate the proposed method on CIFAR-10 [13] with DDIM [32] and  
244 DDPM [10]. We use the model size and OPs (defined in [18]) to evaluate the efficiency of quantized  
245 and full-precision models.

246 For 50-step DDIM sampler, compared with 8-bit PTQ4DM [30], our 4-bit Q-DM achieves a much  
247 larger compression ratio than 8-bit PTQ4DM, but with significant performance improvement (6.89  
248 FID  $\downarrow$  vs. 18.02 FID  $\downarrow$ ). And it is worth noting that the proposed 2-bit model significantly compresses  
249 the DDIM by  $30.9 \times$  on OPs. The proposed method boosts the performance of 4/3/2-bit Baseline  
250 by 3.33, 4.17, and 3.48 in terms of FID score with the same architecture and bit-width, which is  
251 significant on the CIFAR-10 [13] dataset with  $32 \times 32$  generating resolution. For 1000-step DDPM,  
252 the performance of the proposed method outperforms the 4/3/2-bit Baseline by 3.94, 4.14, and 2.58, a  
253 large margin. Also note that the proposed 4/3/2-bit model significantly accelerates the generation by  
254  $7.8 \times$ ,  $15.6 \times$ , and  $30.9 \times$  on OPs. Compared with 8-bit PTQ4DM, ours achieve significantly higher  
255 compression and acceleration rate, while the performance improvement is considerable.

256 Also, our method generates convincing results on ImageNet [14] dataset. As shown in Tab. 2, the  
257 performance of the proposed method with 50-step DDIM significantly outperforms the 4/3/2-bit  
258 Baseline method by 4.76, 4.16, and 4.01. Compared with 8-bit PTQ method, our method achieves  
259 significantly higher compression rate and acceleration rate, but with better performance. For 1000-  
260 step DDPM on ImageNet [14] dataset, the performance of the proposed method outperforms the  
261 4/3/2-bit Baseline by 4.39, 4.54, and 3.71. Also note that our 4-bit Q-DM surpasses the full-precision  
262 50/100-step DDIM and 1000-step DDPM and significantly compresses the noise estimation networks  
263 by  $7.9 \times$ , which demonstrates the effectiveness and efficiency of our Q-DM.

## 264 6 Conclusion

265 In this paper, we present Q-DM, an efficient low-bit quantized diffusion model that offers a high  
266 compression ratio and competitive performance in image generation task. Initially, we analyze the  
267 challenges of the low-bit quantized DM. Our empirical analysis show that distribution oscillation  
268 in activation is the one of the cause of the significant drop in DM quantization. Another challenge  
269 lies in the accumulated quantization error resulted from the multi-step denoising process during  
270 inference. To address these issues, we first develop a timestep-aware quantization (TaQ) method  
271 and a noise-estimating mimicking (NeM) scheme for low-bit quantized DMs, to effectively address  
272 these two challenges. Our work provides a comprehensive analysis and effective solutions for the  
273 crucial issues in low-bit quantized diffusion model, paving the way for the extreme compression and  
274 acceleration of diffusion model.

275 **References**

- 276 [1] Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-dpm: an analytic estimate of the  
277 optimal reverse variance in diffusion probabilistic models. In *Proc. of ICLR*, pages 1–39, 2022.
- 278 [2] Yaohui Cai, Zhewei Yao, Zhen Dong, Amir Gholami, Michael W Mahoney, and Kurt Keutzer.  
279 Zeroq: A novel zero shot quantization framework. In *Proc. of CVPR*, pages 13169–13178,  
280 2020.
- 281 [3] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan.  
282 Wavegrad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*,  
283 2020.
- 284 [4] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm. int8 (): 8-bit matrix  
285 multiplication for transformers at scale. *arXiv preprint arXiv:2208.07339*, 2022.
- 286 [5] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dhar-  
287 mendra S Modha. Learned step size quantization. In *Proc. of ICLR*, pages 1–12, 2019.
- 288 [6] Jun Fang, Ali Shafiee, Hamzah Abdel-Aziz, David Thorsley, Georgios Georgiadis, and Joseph H  
289 Hassoun. Post-training piecewise linear quantization for deep neural networks. In *Proc. of  
290 ECCV*, pages 69–86, 2020.
- 291 [7] Sicheng Gao, Xuhui Liu, Bohan Zeng, Sheng Xu, Yanjing Li, Xiaoyan Luo, Jianzhuang Liu,  
292 Xiantong Zhen, and Baochang Zhang. Implicit diffusion models for continuous super-resolution.  
293 *arXiv preprint arXiv:2303.16491*, 2023.
- 294 [8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
295 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications  
296 of the ACM*, pages 139–144, 2020.
- 297 [9] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
298 Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proc. of  
299 NeurIPS*, pages 1–12, 2017.
- 300 [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Proc.  
301 of NeurIPS*, pages 6840–6851, 2020.
- 302 [11] Alexia Jolicoeur-Martineau, Ke Li, Rémi Piché-Taillefer, Tal Kachman, and Ioannis Mitliagkas.  
303 Gotta go fast when generating data with score-based models. *arXiv preprint arXiv:2105.14080*,  
304 2021.
- 305 [12] Sangil Jung, Changyong Son, Seohyung Lee, Jinwoo Son, Jae-Joon Han, Youngjun Kwak,  
306 Sung Ju Hwang, and Changkyu Choi. Learning to quantize deep networks by optimizing  
307 quantization intervals with task loss. In *Proc. of CVPR*, pages 4350–4359, 2019.
- 308 [13] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
309 Technical report, Citeseer, 2009.
- 310 [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep  
311 convolutional neural networks. In *Proc. of NeurIPS*, pages 1097–1105, 2012.
- 312 [15] Fengfu Li, Bo Zhang, and Bin Liu. Ternary weight networks. *arXiv preprint arXiv:1605.04711*,  
313 2016.
- 314 [16] Yanjing Li, Sheng Xu, Baochang Zhang, Xianbin Cao, Peng Gao, and Guodong Guo. Q-vit:  
315 Accurate and fully quantized low-bit vision transformer. In *Proc. of NeurIPS*, pages 1–12, 2022.
- 316 [17] Yang Lin, Tianyu Zhang, Peiqin Sun, Zheng Li, and Shuchang Zhou. Fq-vit: Fully quantized  
317 vision transformer without retraining. In *Proc. of IJCAI*, pages 1173–1179, 2022.
- 318 [18] Zechun Liu, Zhiqiang Shen, Marios Savvides, and Kwang-Ting Cheng. Reactnet: Towards  
319 precise binary neural network with generalized activation functions. In *Proc. of ECCV*, pages  
320 143–159, 2020.

- 321 [19] Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, and Wen Gao. Post-training  
322 quantization for vision transformer. In *Proc. of NeurIPS*, pages 1–12, 2021.
- 323 [20] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for  
324 improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021.
- 325 [21] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation  
326 with diffusion models. *arXiv preprint arXiv:2103.16091*, 2021.
- 327 [22] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic  
328 models. In *Proc. of ICML*, pages 8162–8171, 2021.
- 329 [23] Chenhao Niu, Yang Song, Jiaming Song, Shengjia Zhao, Aditya Grover, and Stefano Ermon.  
330 Permutation invariant graph generation via score-based generative modeling. In *Proc. of*  
331 *AISTATS*, pages 4474–4484, 2020.
- 332 [24] Haotong Qin, Yifu Ding, Mingyuan Zhang, Qinghua Yan, Aishan Liu, Qingqing Dang, Ziwei  
333 Liu, and Xianglong Liu. Bibert: Accurate fully binarized bert. In *Proc. of ICLR*, pages 1–24,  
334 2022.
- 335 [25] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet  
336 classification using binary convolutional neural networks. In *Proc. of ECCV*, pages 525–542,  
337 2016.
- 338 [26] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad  
339 Norouzi. Image super-resolution via iterative refinement. *TPAMI*, 2022.
- 340 [27] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.  
341 Improved techniques for training gans. In *Proc. of NeurIPS*, pages 1–9, 2016.
- 342 [28] Robin San-Roman, Eliya Nachmani, and Lior Wolf. Noise estimation for generative diffusion  
343 models. *arXiv preprint arXiv:2104.02600*, 2021.
- 344 [29] Hiroshi Sasaki, Chris G Willcocks, and Toby P Breckon. Unit-ddpm: Unpaired image translation  
345 with denoising diffusion probabilistic models. *arXiv preprint arXiv:2104.05358*, 2021.
- 346 [30] Yuzhang Shang, Zhihang Yuan, Bin Xie, Bingzhe Wu, and Yan Yan. Post-training quantization  
347 on diffusion models. *arXiv preprint arXiv:2211.15736*, 2022.
- 348 [31] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu,  
349 Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without  
350 text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- 351 [32] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In  
352 *Proc. of ICLR*, pages 1–20, 2020.
- 353 [33] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data  
354 distribution. In *Proc. of NeurIPS*, pages 1–13, 2019.
- 355 [34] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models.  
356 In *Proc. of NeurIPS*, pages 12438–12448, 2020.
- 357 [35] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and  
358 Ben Poole. Score-based generative modeling through stochastic differential equations. In *Proc.*  
359 *of ICLR*, pages 1–36, 2021.
- 360 [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
361 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. of NeurIPS*, pages 1–11,  
362 2017.
- 363 [37] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks.  
364 In *Proc. of CVPR*, pages 7794–7803, 2018.
- 365 [38] Daniel Watson, Jonathan Ho, Mohammad Norouzi, and William Chan. Learning to efficiently  
366 sample from diffusion probabilistic models. *arXiv preprint arXiv:2106.03802*, 2021.

- 367 [39] Zheng Xie, Zhiqian Wen, Jing Liu, Zhiqiang Liu, Xixian Wu, and Mingkui Tan. Deep  
368 transferring quantization. In *Proc. of ECCV*, pages 625–642, 2020.
- 369 [40] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao,  
370 Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of  
371 methods and applications. *arXiv preprint arXiv:2209.00796*, 2022.
- 372 [41] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. Dorefa-net:  
373 Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint  
374 arXiv:1606.06160*, 2016.
- 375 [42] Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. Trained ternary quantization. In  
376 *Proc. of ICLR*, pages 1–10, 2017.
- 377 [43] Bohan Zhuang, Chunhua Shen, Mingkui Tan, Lingqiao Liu, and Ian Reid. Towards effective  
378 low-bitwidth convolutional neural networks. In *Proc. of CVPR*, pages 7920–7928, 2018.