# Regularized Vector Quantization for Tokenized Image Synthesis

Jiahui Zhang[1]    Fangneng Zhan[2]    Christian Theobalt[2]    Shijian Lu[*1]

[1] Nanyang Technological University    [2] Max Planck Institute for Informatics

## Abstract

*Quantizing images into discrete representations has been a fundamental problem in unified generative modeling. Predominant approaches learn the discrete representation either in a deterministic manner by selecting the best-matching token or in a stochastic manner by sampling from a predicted distribution. However, deterministic quantization suffers from severe codebook collapse and misalignment with inference stage while stochastic quantization suffers from low codebook utilization and perturbed reconstruction objective. This paper presents a regularized vector quantization framework that allows to mitigate above issues effectively by applying regularization from two perspectives. The first is a prior distribution regularization which measures the discrepancy between a prior token distribution and the predicted token distribution to avoid codebook collapse and low codebook utilization. The second is a stochastic mask regularization that introduces stochasticity during quantization to strike a good balance between inference stage misalignment and unperturbed reconstruction objective. In addition, we design a probabilistic contrastive loss which serves as a calibrated metric to further mitigate the perturbed reconstruction objective. Extensive experiments show that the proposed quantization framework outperforms prevailing vector quantization methods consistently across different generative models including auto-regressive models and diffusion models.*

## 1. Introduction

With the prevalence of multi-modal image synthesis [3, 24, 38, 40] and Transformers [32], unifying data modeling regardless of data modalities has attracted increasing interest from the research communities. Aiming for a generic data representation across different data modalities, discrete representation learning [22, 26] plays a significant role in the unified modeling. In particular, vector quantization models (e.g., VQ-VAE [22] and VQ-GAN [8]) emerge as a promising family for learning generic image representations by discretizing images into discrete tokens. With the

---

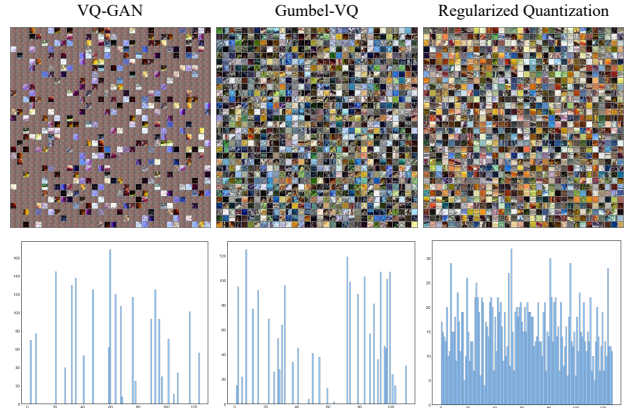*Corresponding author, E-mail: shijian.lu@ntu.edu.sg



Figure 1. Visualization of codebook (first row) and illustration of codebook utilization (second row) on ADE20K dataset [43]. VQ-GAN [8] severely suffers from codebook collapse as most codebook embeddings are invalid values. Gumbel-VQ [2] learns valid values for all codebook embeddings, while only a small number of embeddings are actually used for quantization as illustrated in codebook utilization. As a comparison, the proposed regularized quantization prevents codebook collapse and achieves full codebook utilization. The codebook visualization method is provided in the supplementary file.

tokenized representation, generative models such as auto-regressive model [8, 9] and diffusion model [6, 12] can be applied to accommodate the dependency of the sequential tokens for image generation, which is referred as *tokenized image synthesis* under this context.

Vector quantization models can be broadly grouped into **deterministic** quantization and **stochastic** quantization according to the selection of discrete tokens. Specifically, typical deterministic methods like VQ-GAN [8] directly select the best-matching token via Argmin or Argmax, while stochastic methods like Gumbel-VQ [2] select a token by stochastically sampling from a predicted token distribution. On the other hand, deterministic quantization suffers from codebook collapse [27], a well-known problem where large portion of codebook embeddings are invalid with near-zero values as shown in Fig. 9 (first row). In addition, deterministic quantization is misaligned with the inference stage of generative modeling, where the tokens are usually randomly

sampled instead of selecting the best matching one. Instead, stochastic quantization samples tokens according to a predicted token distribution with Gumbel-Softmax [2, 14], which allows to avoid codebook collapse and mitigate inference misalignment. However, although most codebook embeddings are valid values in stochastic quantization, only a small part is actually utilized for vector quantization as shown in Fig. 9 (second row), which is dubbed as low codebook utilization. Besides, as stochastic methods randomly sample tokens from a distribution, the image reconstructed from the sampled tokens is usually not well aligned with the original image, leading to perturbed reconstruction objective and unauthentic image reconstruction.

In this work, we introduce a regularized quantization framework that allows to prevent above problems effectively via regularization from two perspectives. Specifically, to avoid codebook collapse and low codebook utilization where only a small number of codebook embeddings are valid or used for quantization, we introduce a **prior distribution regularization** by assuming a uniform distribution as the prior for token distribution. As the posterior token distribution can be approximated by the quantization results, we can measure the discrepancy between the prior token distribution and posterior token distribution. By minimizing the discrepancy during training, the quantization process is regularized to use all the codebook embeddings, which prevents the predicted token distribution from collapse into a small number of codebook embeddings.

As deterministic quantization suffers from inference stage misalignment and stochastic quantization suffers from perturbed reconstruction objective, we introduce a **stochastic mask regularization** to strike a good balance between them. Specifically, the stochastic mask regularization randomly masks certain ratio of regions for stochastic quantization, while leaving the unmasked regions for deterministic quantization. This introduces uncertainty for the selection of tokens and results of quantization, which narrows the gap with the inference stage of generative modelling where tokens are selected randomly. We also conduct thorough and comprehensive experiments to analyze the selection of masking ratio for optimal image reconstruction and generation.

On the other hand, with the randomly sampled tokens, the stochastically quantized region will suffer from perturbed reconstruction objective. The perturbed reconstruction objective mainly results from the target for perfect reconstruction of the original image from randomly sampled tokens. Instead of naively enforcing a perfect image reconstruction with L1 loss, we introduce a contrastive loss for elastic image reconstruction, which mitigates the perturbed reconstruction objective significantly. Similar to PatchNCE [23, 42], the contrastive loss treats the patch at the same spatial location as positive pairs and others as neg-

ative pairs. By pushing the positive pairs closer and pulling negative pairs away, the elastic image reconstruction can be achieved. Another issue with the randomly sampled tokens is that they tend to introduce perturbation of different scales in the reconstruction objective, We thus introduce a Probabilistic Contrastive Loss (PCL) that adjusts the pulling force of different regions according to the discrepancy between the sampled token embedding and the best-matching token embedding.

The contributions of this work can be summarized in three aspects. First, we present a regularized quantization framework that introduces a prior distribution regularization to prevent codebook collapse and low codebook utilization. Second, we propose a stochastic mask regularization which mitigates the misalignment with the inference stage of generative modelling. Third, we design a probabilistic contrastive loss that achieves elastic image reconstruction and mitigates the perturbed objective adaptively for different regions with stochastic quantization.

## 2. Related Work

### 2.1. Vector Quantization

As introduced in Oord *et al.* [22], vector quantization aims to represent the data with entries of a learnt codebook (i.e., tokens), which achieves discrete and compressed representation. According to the selection mechanism of discrete tokens, the quantization methods can be grouped into deterministic quantization and stochastic quantization.

**Deterministic Quantization.** With a predicted token distribution or probability, deterministic quantization aims to select the best-matching token through Argmin or Argmax. Typically, VQ-VAE [22] proposes to quantize the encoded feature into discrete token by looking up the nearest neighbour (i.e., Argmin) entry in a learned codebook. As the operation of Argmax is not differentiable, there is no real gradient defined for the encoder. VQ-VAE adopts Straight Through Estimation (STE) [4] by copying the gradient from the decoder to the encoder. During training, the encoder, decoder and codebook are optimized driven by a reconstruction loss and a codebook embedding loss [22]. Following this line of deterministic quantization, strenuous effort has been made to improve the quantization performance, including Exponential Moving Averages (EMA) [22, 25] for stable updating of codebook, multi-scale hierarchical organization [25] for higher synthetic coherence and fidelity, adversarial loss and perceptual loss [8] for improved perceptual quality, integrated quantization [39] for condition generation, translation-equivariant quantization with orthogonal codebook embeddings [28], Transformer structure for quantization [36].

**Stochastic Quantization.** Instead of naively selecting the best-matching token, stochastic quantization aims

to sample token from a predicted token distribution. As the sampling operation is not differentiable, certain reparameterization trick (e.g., Gumbel-Softmax [14]) should be applied for gradient backpropagation. For instance, VQ-Wave2Vec [2] introduces stochastic quantization with Gumbel-Softmax reparameterization trick to learn the discrete representation of audios. Similarly, DALL-E [24] leverages Gumbel-Softmax to represent images with tokens, followed by a Transformer to auto-regressively model the dependency between tokens for diverse image synthesis. Inspired by the connection between Exponential Moving Averages (EMA) [22] and Expectation Maximization (EM) algorithm [20], Roy *et al.* [27] introduce soft EM algorithm for quantization by performing Monte-Carlo Expectation Maximization [34] on the probability distribution of discrete latent variables. The sampling from a Gumbel-Softmax distribution exactly approximates the sampling from the token distribution as proved in [19].

## 2.2. Tokenized Image Synthesis

With a tokenized representation of images, generative models can be applied to the discrete tokens for image synthesis [8, 37, 39]. For example, PixelRNN and PixelCNN [31] employ LSTM [13] and masked convolutions to model inter-dependencies of pixels autoregressively. By learning a compressed tokenized representation, VQ-VAE [22] achieves compelling generation quality with PixelCNN. Recently, Transformer [32] emerges as a powerful paradigm for sequence modeling. Chen *et al.* [5] leverage Transformer to model the sequence dependency of image pixels. DALL-E [24] designs a discrete VAE for learning tokenized representation and models the image tokens with a Transformer to achieve text-to-image generation. Esser *et al.* [8] propose a VQ-GAN to learn a rich discrete representation and utilize the Transformer to efficiently model token distributions for high-resolution images synthesis.

In addition to auto-regressive models, diffusion models [6, 12, 30] can also be applied to model discrete tokens. For instance, D3PMs [1] adopts discrete diffusion to estimate the density of image pixels and achieves low-resolution image synthesis. With the learned discrete tokens in VQ-VAE, VQ-Diffusion [10] achieves compelling text-to-image generation performance via a discrete diffusion process.

## 3. Method

As illustrated in Fig. 2, our regularized quantization framework combines deterministic quantization and stochastic quantization, which consists of an encoder $E$, a decoder $G$, and a codebook $\mathcal{Z} = \{z_n\}_{n=1}^N \in \mathbb{R}^{N \times d}$, where $N$ is the size of the codebook and $d$ is the dimension of embeddings. Given an input image $X$, the encoder is employed to produce a spatial collection of token distributions

$x_i \in \mathbb{R}^N, i \in [1, H \times W]$, where $H \times W$ is the size of spatial vectors. Then, each encoded vector is mapped into a discrete token according to the predicted token distribution, which yields tokenized representation (i.e., indices of codebook embeddings). The codebook embeddings associated with the indices are finally fed into the decoder to reconstruct the input image.

With a trained vector quantization framework, we can represent images in terms of the codebook indices (i.e., tokens). With the discrete tokens of images, generative models such as auto-regressive model [9] and diffusion model [12] can be applied to build the dependency between tokens. At inference stage of generative modeling, a sequence of tokens can be sampled for image synthesis. By mapping the sequence of tokens back to their corresponding codebook embeddings, an image can be readily generated by feeding the embeddings into the decoder.

### 3.1. Prior Distribution Regularization

Prevailing vector quantization models usually severely suffer from codebook collapse or low codebook utilization, where only a small number of codebook embeddings are valid or used for quantization. We thus propose a prior distribution regularization to regularize the vector quantization process. Specifically, we assume a prior distribution for the tokens used for quantization. Ideally, the prior distribution is expected to be a uniform discrete distribution denoted by $P_{prior} = [1/N, 1/N, \cdots, 1/N], P_{prior} \in \mathbb{R}^N$, which means all codebook embeddings can be used uniformly and their corresponding information capacity can be maximized according to the principle of maximum entropy. During quantization, as image features of size $H \times W$ are mapped to corresponding tokens, the predicted quantization result of each feature can be represented by a one-hot vector $p_i, i \in [1, H * W]$. Thus, the posterior token distribution $P_{post}$ can be approximated by the average of all one-hot vector as $P_{post} = \sum_{i=1}^{H \times W} p_i / (H \times W) = [p_1, p_2, \cdots, p_N]$. Then, the discrepancy between the prior token distribution and predicted token distribution can be measured by KL divergence as below:

$$\mathcal{L}_{kl} = KL(P_{post}, P_{prior}) = -\sum_n^N p_n \log \frac{1/N}{p_n} \quad (1)$$

By minimizing the KL divergence $\mathcal{L}_{kl}$, the vector quantization can be effectively regularized to avoid codebook collapse and low codebook utilization.

### 3.2. Stochastic Mask Regularization

On the other hand, the deterministic quantization which selects the most probable tokens will lead to misalignment with the inference stage of generative modelling, where tokens are sampled randomly according to the predicted dis-
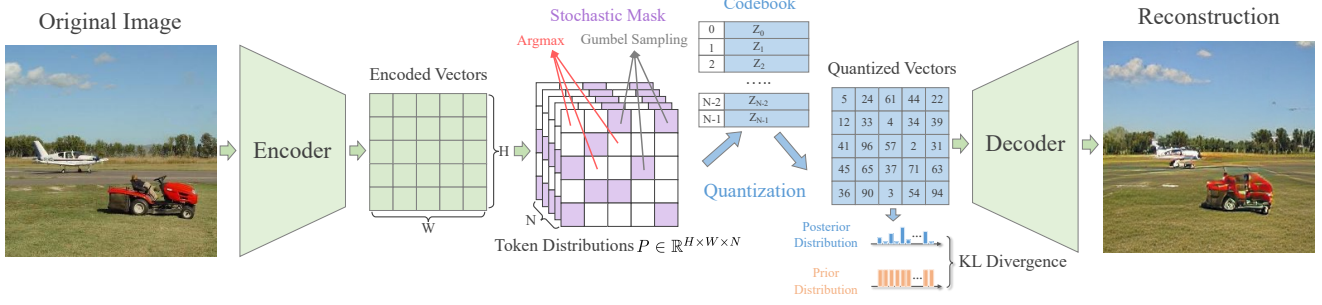
Figure 2. The framework of the proposed regularized quantization. A stochastic mask (indicated by purple regions) is applied to the *Predicted Token Distributions* to specify the region for stochastic sampling. Then the *Encoded Vectors* can be represented by the selected codebook embeddings, which produces the *Quantized Vectors* for image *Reconstruction*. A KL divergence is measured between the posterior token distribution and the prior token distribution to avoid codebook collapse and low codebook utilization.
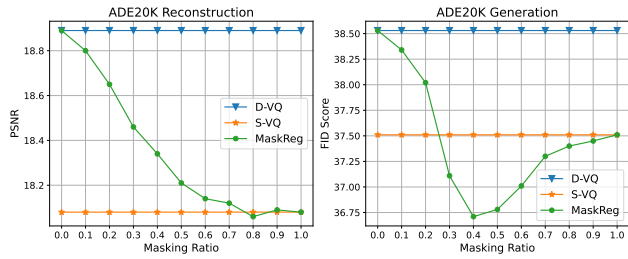


Figure 3. The effect of different masking ratios on image reconstruction and image generation on ADE20K dataset. VQ-GANs with different quantization methods are used for image tokenization and auto-regressive model is used for image synthesis. D-VQ, S-VQ, and MaskReg denote deterministic vector quantization, stochastic vector quantization, and the proposed stochastic mask regularization, respectively. The quality of image reconstruction is evaluated by PSNR, the quality of image generation is evaluated by FID. Note good PSNR score doesn't indicate good generation quality.
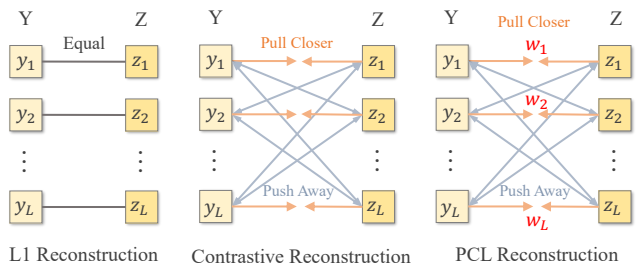


Figure 4. Comparisons of L1 loss, vanilla contrastive loss, and our proposed probabilistic contrastive loss (PCL) for image reconstruction. PCL introduces adaptive weights $\{w_i\}_{i=1}^{N}$ to positive pairs for better representation learning and elastic image reconstruction.

tribution. Instead, stochastic quantization enables to introduce stochasticity during quantization which helps to mitigate the misalignment with inference stage. Nevertheless, stochastic quantization tends to incur perturbed reconstruction objective as the sampled tokens may not match the original image. As a result, the generation quality (FID) of stochastic quantization presents marginal improvement compared with deterministic quantization as shown in Fig. 3. To strike a good balance between unperturbed reconstruction objective and inference stage misalignment, we design a stochastic mask regularization which combines the deterministic quantization and stochastic quantization by applying a stochastic mask to coordinate the image regions for stochastic quantization and deterministic quantization as shown in Fig. 2.

Specifically, with the predicted token probability $P \in \mathbb{R}^{H \times W \times N}$ for all encoded vectors, we randomly set a mask

$M \in \mathbb{R}^{H \times W}$ with '1' to indicate the regions for token sampling with Gumbel-softmax and '0' to indicate the regions for selecting the best-matching tokens with Argmax. Denoting the vectors quantized via Argmax and Gumbel sampling as $X_{argmax}$ and $X_{gumbel} \in \mathbb{R}^{H \times W \times N}$, respectively, the reconstruction objective can be formulated as below:

$$\mathcal{L}_{rec} = \|X - G(X_{argmax} * (1 - M) + X_{gumbel} * M)\|_1,$$

where $G$ denotes the decoder. As shown in Fig. 3, comprehensive experiments are conducted to analyze the effect with different masking ratios, and a masking ratio of 40% is proved to yield the best image reconstruction & generation quality (i.e., best FID). The effect of different masking ratios on other datasets is provided in the supplementary material. As both Argmax and Gumbel sampling operations are non-differentiable, we apply reparameterization trick by replacing Argmax operation with Softmax and Gumbel with Gumbel-Softmax in gradient back-propagation. The pseudo code of the forward & backward propagation of the proposed regularized quantization is given in Algorithm 1.

**Algorithm 1** Pseudo-code of forward & backward propagation in vector quantization with the proposed stochastic mask regularization.

---

**Deterministic Quantization Region.**
**Input:** encoded vector $x_{ij}$, token distribution $P_{ij} = [p_1, p_2, \cdots, p_N]$, codebook $\mathcal{Z}$.
    **Forward propagation:**
    1. index = Argmax($P_{ij}$)
    2. index_hard = One_Hot(index)
    3. quantized $\hat{x}_{ij}$ = Matmul(index_hard, $\mathcal{Z}$)
    **Backward propagation:**
    1. index_soft = Softmax($P_{ij}$)
    2. quantized $\hat{x}_{ij}$ = Matmul(index_soft, $\mathcal{Z}$)
**Output:** quantized $\hat{x}_{ij}$.

**Stochastic Quantization Region.**
**Input:** encoded vector $x_{ij}$, token distribution $P_{ij} = [p_1, p_2, \cdots, p_N]$, codebook $\mathcal{Z}$, gumbels $\sim$ Gumbel(0, 1).
    **Forward propagation:**
    1. index = Argmax($P_{ij}$+gumbels)
    2. index_hard = One_Hot(index)
    3. quantized $\hat{x}_{ij}$ = Matmul(index_hard, $\mathcal{Z}$)
    **Backward propagation:**
    1. index_soft = Softmax($P_{ij}$+gumbels)
    2. quantized $\hat{x}_{ij}$ = Matmul(index_soft, $\mathcal{Z}$)
**Output:** quantized $\hat{x}_{ij}$.

---

### 3.3. Probabilistic Contrastive Loss

The stochastic mask regularization mitigates the misalignment with the inference stage. However, for the image region with stochastic quantization, the model training still suffers from perturbed reconstruction objective caused by the randomly sampled tokens. Thus, we propose a probabilistic contrastive loss (PCL) to mitigate the perturbed objective in the region with stochastic quantization. As the perturbed objective results from the target for perfect reconstruction of original images with L1 loss, the proposed PCL achieves *elastic image reconstruction* [1] in the stochastic quantization region through contrastive learning.

Instead of forcing perfect image reconstruction, contrastive learning aims to maximize the mutual information between corresponding images by pulling selected positive pairs closer and pushing negative pairs away as shown in Fig. 4. Following the Noise Contrastive Estimation framework [21] in PatchNCE [23,41], image features in the same spatial location of the original and reconstructed image are regarded as positive pairs and others are negative pairs in PCL. Thus, vanilla contrastive loss $\mathcal{L}_{cl}$ for image reconstruction can be formulated as:

$$\mathcal{L}_{cl} = -\frac{1}{L} \sum_{i=1}^{L} \log \frac{e^{y_i \cdot z_i / \tau}}{e^{y_i \cdot z_i / \tau} + \sum_{\substack{j=1 \\ j \neq i}}^{L} e^{y_i \cdot z_j / \tau}}, \quad (2)$$

where $Y = [y_1, y_2, \cdots, y_L]$ and $Z = [z_1, z_2, \cdots, z_L]$ are

---

[1]'Elastic' means relative / contrastive approximation of the original image instead of perfect reconstruction.

extracted feature patches from the original image and reconstructed image respectively, $\tau$ is the temperature parameter, $L$ is the number of image features. As proved in [7, 8], perceptual loss [15] helps to keep good perceptual quality for image reconstruction. We thus employ pretrained VGG-19 network [29] to extract multi-layer image features ($relu1\_2, relu2\_2, relu3\_3, relu4\_3, relu5\_3$) from original and reconstructed images to construct contrastive learning pairs. Note, the contrastive loss is used alongside the perceptual loss during training.

**Probabilistic Contrast.** With the stochasticity in Gumbel sampling, the sampled tokens tend to present varying discrepancies with the best-matching one as selected by Argmax. Intuitively, a sampled token with larger discrepancy with the best-matching one will yield severer objective perturbation. Thus, the pulling force between the original and reconstructed images should be adaptive with respect to the perturbation magnitude for optimal contrastive learning. We introduce the Probabilistic Contrastive Loss (PCL) which employs **weighting parameters** $\{w_i\}_{i=1}^{L}$ to adjust the pulling force of different features according to the token sampling results (i.e., perturbation magnitude) as shown in Fig. 4. The weighting parameter $w_i$ is produced by computing the Euclidean distance between the randomly sampled embedding (denoted by $z_s$) and the best-matching embedding (denoted by $z_q$): $w_i = \|z_s - z_q\|_2^2$. Then, the probabilistic contrastive loss $\mathcal{L}_{pcl}$ can be formulated by adjusting the pulling force of positive pairs with the normalized weighting parameters $\{w_i'\}_{i=1}^{L}, s.t. \sum_{i=1}^{N} w_i' = 1$ as below:

$$\mathcal{L}_{pcl} = -\sum_{i=1}^{L} \log \frac{w_i' \cdot e^{y_i \cdot z_i / \tau}}{w_i' \cdot e^{y_i \cdot z_i / \tau} + \frac{1}{L} \sum_{\substack{j=1 \\ j \neq i}}^{L} e^{y_i \cdot z_j / \tau}}. \quad (3)$$

Note, we balance the negative term with $1/L$, otherwise the negative term will be too large compared with the vanilla contrastive loss.

## 4. Experiments

### 4.1. Experimental Settings

**Datasets** We benchmark our method over multiple public datasets, including ADE20K [43] and CelebA-HQ [17] for semantic image synthesis, CUB-200 [35] and MS-COCO [16] for text-to-image synthesis.

**Evaluation Metrics.** We evaluate the vector quantization models by assessing their image reconstruction and image generation performance. The image reconstruction performance is evaluated with several widely adopted evaluation metrics. Specifically, Fréchet Inception Score (FID) [11] is employed to evaluate the quality (perceptual similarity) of reconstructed images and generated images; Peak Signal-to-noise Ratio (PSNR) is employed to measure

Table 1. Semantic image synthesis with **auto-regressive** models on ADE20K and CelebA-HQ, and text-to-image synthesis with **diffusion** models on CUB-200 and MS-COCO. [R] and [G] denote the results of reconstructed images, generated images with auto-regressive models or diffusion models.

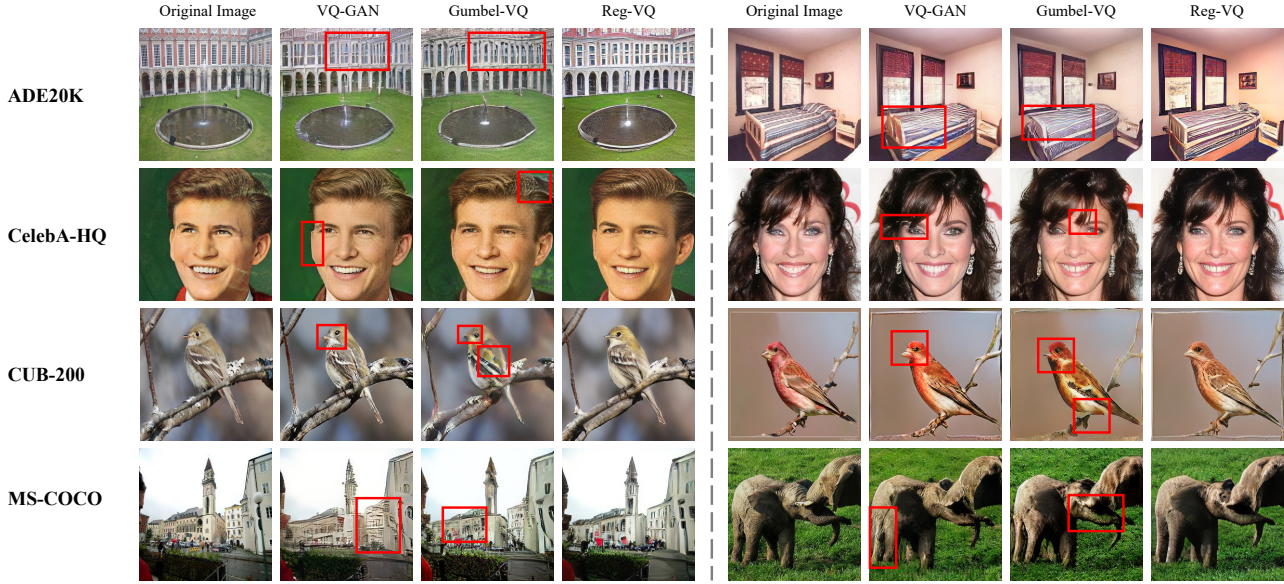| Models | ADE20K [43] (Semantic) | | | CelebA-HQ [17] (Semantic) | | | CUB-200 [35] (Text) | | | MS-COCO [16] (Text) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FID[R]↓ | PSNR[R]↑ | FID[G]↓ | FID[R]↓ | PSNR[R]↑ | FID[G]↓ | FID[R]↓ | PSNR[R]↑ | FID[G]↓ | FID[R]↓ | PSNR[R]↑ | FID[G]↓ |
| **VQ-VAE** [33] | 49.21 | **19.95** | 60.29 | 28.38 | **23.39** | 39.57 | 20.89 | **21.07** | 26.32 | 32.48 | **19.12** | 38.84 |
| **VQ-GAN** [8] | 28.17 | 18.89 | 38.53 | 12.74 | 22.44 | 17.42 | 13.49 | 20.88 | 17.43 | 18.58 | 18.86 | 23.75 |
| **Gumbel-VQ** [2] | 26.42 | 18.08 | 37.51 | 12.03 | 20.92 | 16.78 | 13.25 | 20.06 | 16.93 | 16.97 | 18.43 | 22.21 |
| **Reg-VQ** | **23.69** | 18.44 | **34.47** | **10.09** | 22.05 | **15.34** | **10.84** | 20.39 | **14.14** | **13.76** | 18.64 | **19.91** |



Figure 5. Reconstruction of images from four public datasets with different quantization methods: The red-color boxes highlight reconstruction artifacts.

the accuracy (pixel-level similarity) of image reconstruction.

**Implementation Details.** Following VQ-GAN [8], the default feature size $H \times W$ and codebook size $N$ for all methods are set as $16 \times 16$ and 1024, respectively [2]. An image size of $256 \times 256$ is adopted for both image reconstruction and image generation. A masking ratio of 40% is adopted in our experiments by default. More details of the experiments (e.g., hyper-parameters, network architectures) are provided in the supplemental material.

### 4.2. Quantitative Evaluation

Vector quantization performance can be evaluated by assessing their image reconstruction and image generation performance. For image reconstruction, quantized image tokens are fed into the decoders of quantization models to recover the original images. For image generation, auto-regressive model [8] [3] is employed for semantic image synthesis on ADE20K and CelebA-HQ; diffusion model [10] [4] is employed for text-to-image synthesis on CUB-200 and MS-COCO. Details of the used auto-regressive model and diffusion model are provided in the supplementary material.

Table 1 shows the image reconstruction & generation results on ADE20K & CelebA-HQ and CUB-200 & MS-COCO. VQ-GAN [8] is a deterministic quantization method, while Gumbel-VQ is a stochastic variant of VQ-GAN by employing Gumbel-Softmax for token sampling as in DALL-E [24]. It can be observed that the proposed regularized quantization (Reg-VQ) outperforms all compared methods in terms of the reconstruction quality and generation quality as evaluated by FID[R] and FID[G], respectively. VQ-GAN achieves relatively high reconstruction accuracy as evaluated by PSNR[R], as it naively selects

---

[2]Thus, the results reported in our paper are different from that in VQ-Diffusion whose default feature size and codebook size are $32 \times 32$ and 8192, respectively.

[3]https://github.com/CompVis/taming-transformers
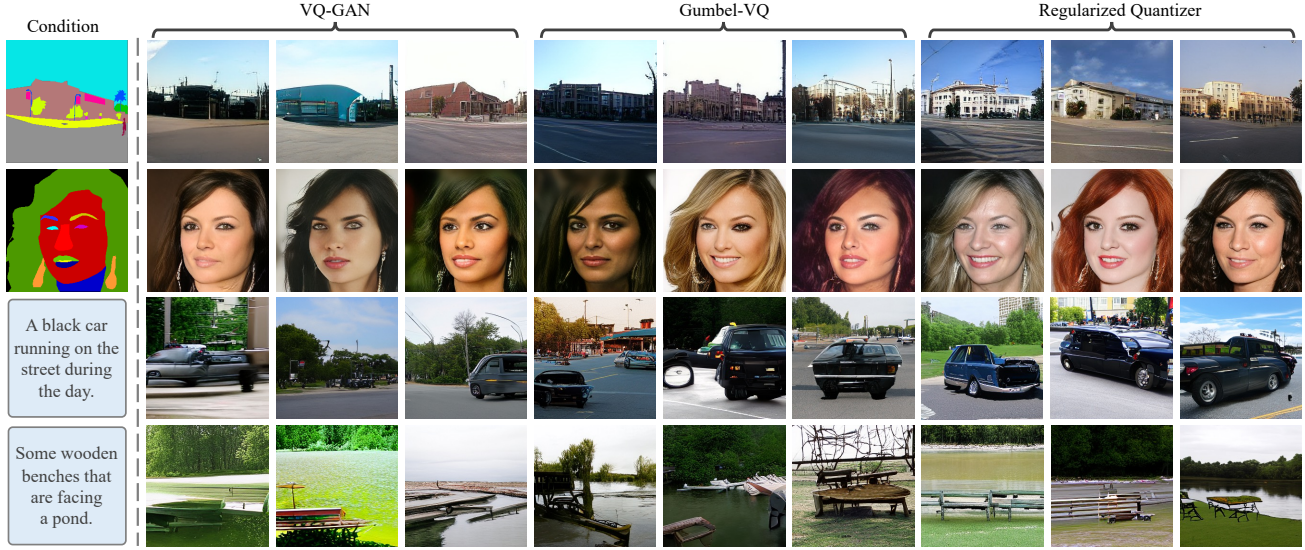[4]https://github.com/cientgu/VQ-Diffusion

Figure 6. Semantic image synthesis and text-to-image synthesis. Three synthesis samples are shown for each condition under each quantization method.

the best-matching token for reconstruction. However, the high reconstruction accuracy doesn't contribute to the image generation performance as reconstructing the original images is not the objective of image generation. Instead, we can observe that it is the high **reconstruction quality** (as evaluated by FID[R]) that indicates a high generation performance (as evaluated by FID[G]), which can be further validated in the ensuing qualitative evaluation.

## 4.3. Qualitative Evaluation

We qualitatively compare the image reconstruction and generation performance of different methods as shown in Fig. 5 and Fig. 6. Regularized quantization (Reg-VQ) achieves the best reconstruction quality although its reconstructed images are not exactly aligned with original images in terms of detailed textures. As suffering from codebook collapse or low codebook utilization, both VQ-GAN and Gumbel-VQ present inferior reconstruction quality. Regularized quantization also achieves superior synthesis quality on various image generation tasks (semantic image synthesis and text-to-image synthesis) and generative models (auto-regressive model and diffusion model) as illustrated in Fig. 6.

## 4.4. Ablation Study

We conduct extensive ablation studies to evaluate the regularized quantization as shown in Table 2. The deterministic quantization method VQ-GAN [8] serves as the baseline model. With the including of Prior distribution regularization (denoted by PriorReg), the image reconstruction quality & accuracy and generation quality are all improved substantially as evaluated by FID[R] & PSNR[R]

Table 2. Ablation study on semantic image synthesis with auto-regressive models. VQ-GAN [8] serves as the baseline model. 'PriorReg', 'MaskReg', 'CL', 'PCL' denote the prior distribution regularization, stochastic mask regularization, vanilla contrastive loss, and our probabilistic contrastive loss, respectively. The row in grey denotes the result of the standard regularized quantization.

| Models | ADE20K [43] (Semantic) | | |
|---|---|---|---|
| | FID[R]↓ | PSNR[R]↑ | FID[G]↓ |
| **Baseline [8]** | 28.17 | 18.89 | 38.53 |
| **+PriorReg** | 25.92 | **18.98** | 36.57 |
| **+PriorReg+MaskReg** | 25.11 | 18.56 | 35.03 |
| **+PriorReg+MaskReg+CL** | 24.21 | 18.49 | 34.91 |
| **+PriorReg+MaskReg+PCL** | **23.69** | 18.44 | **34.47** |

and FID[G]. The including of stochastic mask regularization (denoted by MaskReg) further improves the reconstruction and generation quality but degrades the reconstruction accuracy as shown in **+PriorReg+MaskReg**. To mitigate the perturbed objective, contrastive loss (CL) is included and brings certain performance gains for image reconstruction and generation quality. As a comparison, including the designed probabilistic contrastive loss (PCL) improves the reconstruction and generation performance by a larger margin.

We also study the effect of varying codebook size on ADE20K. As shown in Fig. 7, the vanilla VQ-GAN with deterministic quantization suffers from severer codebook collapse with increasing of codebook size, while the proposed Regularized quantization still achieves high code-
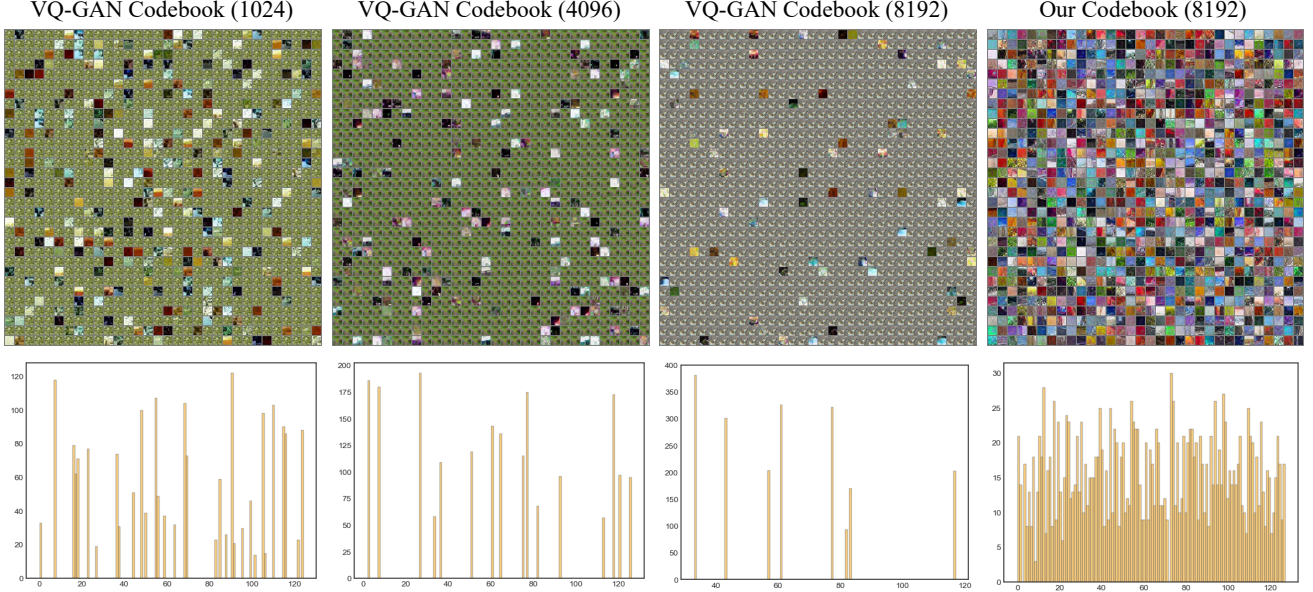
Figure 7. Visualization of codebook with different sizes (N=1024, 4096, 8192) on ADE20K. The first 1024 codebook embeddings are illustrated for all models. Vanilla VQ-GAN (with deterministic quantization) suffers from severe codebook collapse with the increase of codebook size, while the regularized quantization achieves high codebook utilization consistently. The visualization of codebook on other datasets is provided in the supplementary material.
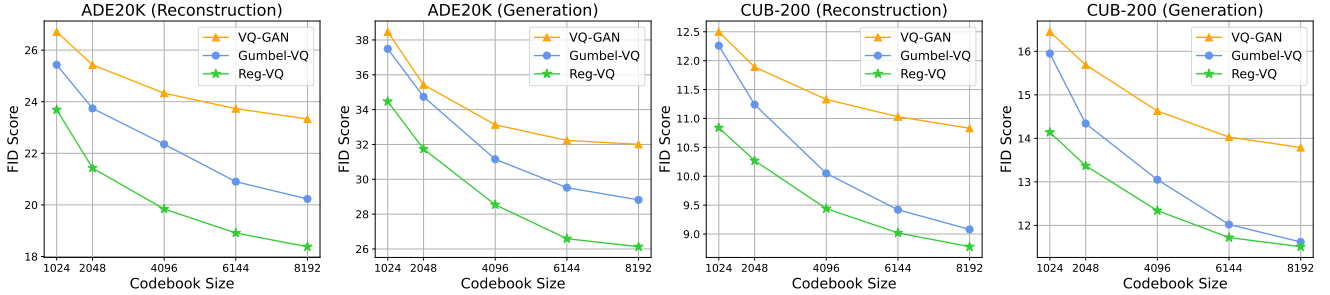


Figure 8. The effect of different codebook size on image reconstruction and image generation. Auto-regressive model is employed for semantic image generation on ADE20K dataset and diffusion model is employed for text-to-image generation on CUB-200.

book utilization for large codebook size. Fig. 8 quantitatively illustrates the reconstruction and generation performance with the different codebook sizes. We can observe that the performance of regularized quantization (Reg-VQ) improves clearly with the increasing of codebook size, while the performance of VQ-GAN tends to be capped at a lower level.

## 5. Conclusions

This paper presents a regularized quantization framework which achieves superior image quantization performance. A prior distribution regularization is proposed to prevent the codebook collapse and low codebook utilization. A stochastic mask regularization is designed to balance the inference stage misalignment and unperturbed re-

construction objective, and the masking ratio is analyzed comprehensively to yield the best performance. To mitigate the perturbed reconstruction objective, a probabilistic contrastive loss is proposed to serve as a calibrated metric for elastic image reconstruction. Quantitative and qualitative experiments show that regularized quantization enables to synthesize high-fidelity images for various generation tasks.

## 6. Acknowledgments

# References

[1] Jacob Austin, Daniel Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34, 2021. 3

[2] Alexei Baevski, Steffen Schneider, and Michael Auli. vq-wav2vec: Self-supervised learning of discrete speech representations. *arXiv preprint arXiv:1910.05453*, 2019. 1, 2, 3, 6

[3] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(2):423–443, 2018. 1

[4] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. 2

[5] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *International Conference on Machine Learning*, pages 1691–1703. PMLR, 2020. 3

[6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34, 2021. 1, 3

[7] Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen, Fang Wen, and Nenghai Yu. Peco: Perceptual codebook for bert pre-training of vision transformers. *arXiv preprint arXiv:2111.12710*, 2021. 5

[8] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 1, 2, 3, 5, 6, 7, 11

[9] Karol Gregor, Ivo Danihelka, Andriy Mnih, Charles Blundell, and Daan Wierstra. Deep autoregressive networks. In *International Conference on Machine Learning*, pages 1242–1250. PMLR, 2014. 1, 3

[10] 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. *arXiv preprint arXiv:2111.14822*, 2021. 3, 6, 11

[11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, pages 6626–6637, 2017. 5

[12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 1, 3

[13] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 3

[14] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. 2, 3

[15] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 694–711. Springer, 2016. 5

[16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2014. 5, 6, 11

[17] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3730–3738, 2015. 5, 6, 11

[18] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 11

[19] Chris J Maddison, Daniel Tarlow, and Tom Minka. A* sampling. In *NIPS*, 2014. 3

[20] Todd K Moon. The expectation-maximization algorithm. *IEEE Signal processing magazine*, 13(6):47–60, 1996. 3

[21] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 5

[22] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. *arXiv preprint arXiv:1711.00937*, 2017. 1, 2, 3

[23] Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation. In *European Conference on Computer Vision*, pages 319–345. Springer, 2020. 2, 5

[24] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 1, 3, 6

[25] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019. 2

[26] Jason Tyler Rolfe. Discrete variational autoencoders. *arXiv preprint arXiv:1609.02200*, 2016. 1

[27] Aurko Roy, Ashish Vaswani, Arvind Neelakantan, and Niki Parmar. Theory and experiments on vector quantized autoencoders. *arXiv preprint arXiv:1805.11063*, 2018. 1, 3

[28] Woncheol Shin, Gyubok Lee, Jiyoung Lee, Joonseok Lee, and Edward Choi. Translation-equivariant image quantizer for bi-directional image-text generation. *arXiv preprint arXiv:2112.00384*, 2021. 2

[29] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5

[30] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 3

[31] Aaron Van Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International Conference on Machine Learning*, pages 1747–1756. PMLR, 2016. 3

[32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 1, 3

[33] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. 6

[34] Greg CG Wei and Martin A Tanner. A monte carlo implementation of the em algorithm and the poor man's data augmentation algorithms. *Journal of the American statistical Association*, 85(411):699–704, 1990. 3

[35] Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona. Caltech-ucsd birds 200. *California Institute of Technology*, 2010. 5, 6, 11

[36] Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved vqgan. *arXiv preprint arXiv:2110.04627*, 2021. 2

[37] Yingchen Yu, Fangneng Zhan, Rongliang Wu, Jianxiong Pan, Kaiwen Cui, Shijian Lu, Feiying Ma, Xuansong Xie, and Chunyan Miao. Diverse image inpainting with bidirectional and autoregressive transformers. *arXiv preprint arXiv:2104.12335*, 2021. 3

[38] Yingchen Yu, Fangneng Zhan, Rongliang Wu, Jiahui Zhang, Shijian Lu, Miaomiao Cui, Xuansong Xie, Xian-Sheng Hua, and Chunyan Miao. Towards counterfactual image manipulation via clip. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 3637–3645, 2022. 1

[39] Fangneng Zhan, Yingchen Yu, Rongliang Wu, Jiahui Zhang, Kaiwen Cui, Changgong Zhang, and Shijian Lu. Autoregressive image synthesis with integrated quantization. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVI*, pages 110–127. Springer, 2022. 2, 3, 11

[40] Fangneng Zhan, Yingchen Yu, Rongliang Wu, Jiahui Zhang, and Shijian Lu. Multimodal image synthesis and editing: A survey. *arXiv preprint arXiv:2112.13592*, 2021. 1

[41] Fangneng Zhan, Yingchen Yu, Rongliang Wu, Jiahui Zhang, Shijian Lu, and Changgong Zhang. Marginal contrastive correspondence for guided image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10663–10672, 2022. 5

[42] Fangneng Zhan, Jiahui Zhang, Yingchen Yu, Rongliang Wu, and Shijian Lu. Modulated contrast for versatile image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18280–18290, 2022. 2

[43] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 633–641, 2017. 1, 5, 6, 7

## A. Codebook Visualization

The codebook is visualized by directly feeding the embeddings in the codebook into the decoder $G$ to generate codebook image patches. Note the decoder $G$ is a fully convolution network and thus the input features to the decoder can have any sizes. For all models, the first 1024 codebook embeddings (i.e., indices 0-1023) are visualized to form a image with $32 \times 32$ patches.

We visualize the codebook of VQ-GAN and our regularized quantizer on CelebA-HQ [17], CUB-200 [35], and MS-COCO [16] as shown in Fig. 9.

## B. Experiment Details

Both vector quantization methods and generative models are optimized via AdamW [18] solver ($\beta_1 = 0.9$ and $\beta_2 = 0.95$) with a learning rate of $1.5e$-4. All experiments are conducted on 4 Tesla V100 GPUs with a batch size of 40.

For vector quantization, the feature size $H \times W$ is set as $16 \times 16$ by default; the default codebook size $N$ and embedding dimension are set as 1024 and 256. The encoder and decoder in regularized quantizer follow the default structure of VQ-GAN [8]. The temperature parameter in Gumbel-Softmax is set as 0.9 by default. The training epochs for vector quantization on ADE20K, CelebA-HQ, CUB-200, and MS-COCO are 100, 60, 300, 50, respectively, for all models.

For auto-regressive modeling, the Transformer used in Esser *et al.* [8,39] [5] is selected as the code base with the default setting. Specifically, the vocabulary size, embedding number and input sequence length are 1024, 1024 and 512, respectively; the numbers of transformer blocks and attention head are 24 and 16, respectively. For the task of semantic image synthesis, the *semantic ids* of semantic maps are directly used as the conditional tokens for the Transformer. The training epoch for the auto-regressive model is 50 for all datasets.

For the diffusion modeling, Denoising Diffusion Probabilistic Model (DDPM) in VQ-Diffusion [10] [6] is selected as the code base. Specifically, VQ-Diffusion-B (Base) which has 19 transformer blocks with dimension of 1024 is employed to estimate the token distribution. The training epoch for the diffusion model is 100 for all datasets.

## C. Limitations and Future Work

Current quantization models train the encoder and decoder with the same learning objective. However, for tokenized image synthesis, the encoder and decoder in quantization models actually have different objectives: the encoder aims to learn accurate discrete representation, while the decoder aims to generate realistic images. Thus, training them with the same objective tends to be sub-optimal and will constrain the quantization and generation performance. In the future, we will explore to design separated learning objective for the encoder and decoder for optimal quantization and generation performance. Besides, we can also explore the performance with different prior distribution, e.g., Gaussian distribution.

## D. Ethical Considerations

The proposed quantization method aims to synthesize realistic images. There would be negative impacts if it is combined with generation methods for certain illegal purpose such as image forgery.

## E. More Qualitative Results

We provide more image translation results including Figs. 10, 11 for semantic image synthesis, and Figs. 12, 13 for text-to-image synthesis.
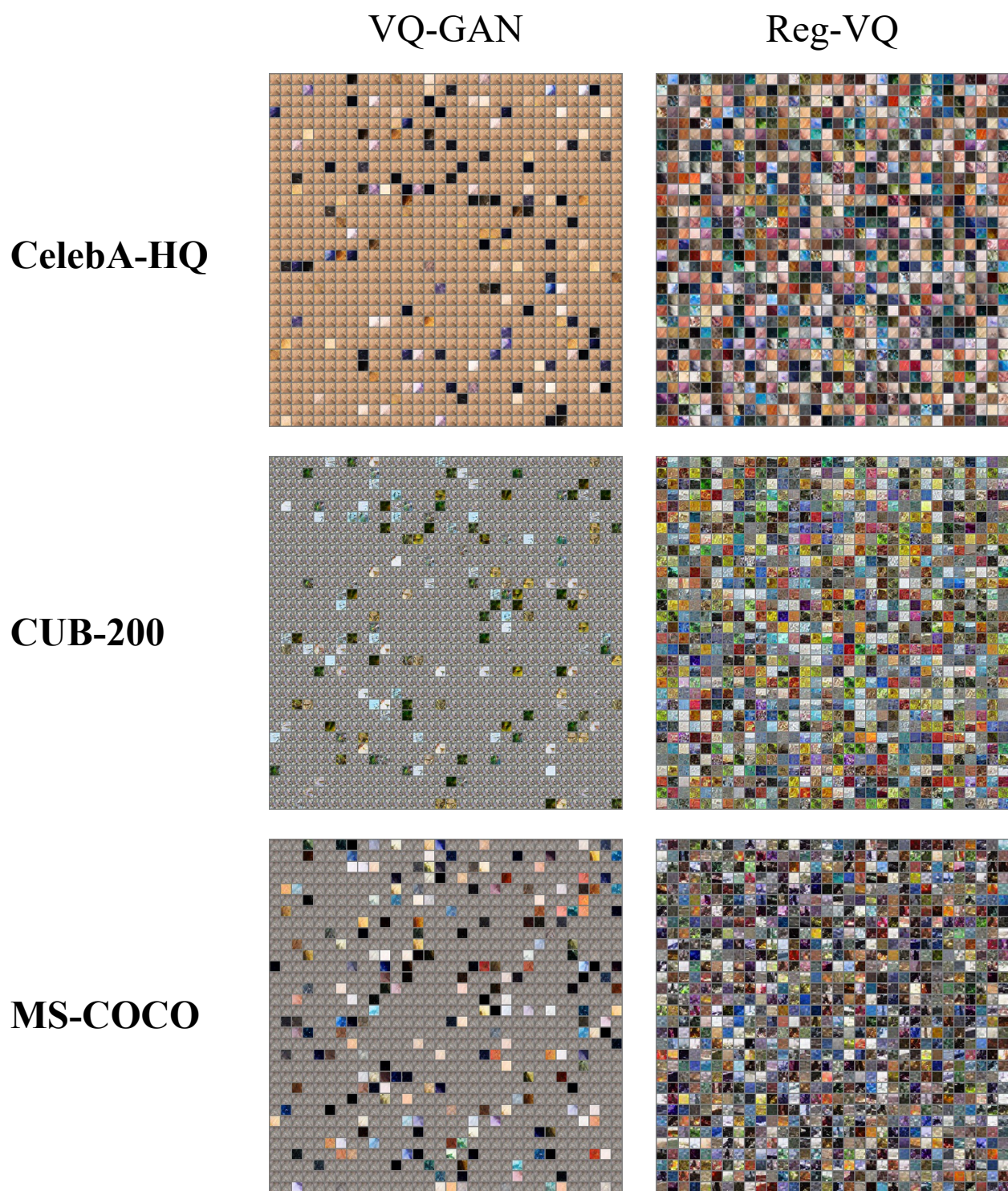
---

[5]https://github.com/CompVis/taming-transformers
[6]https://github.com/cientgu/VQ-Diffusion

Figure 9. The visualization of codebook of VQ-GAN and regularized quantizer on CelebA-HQ, CUB-200, and MS-COCO.

| Input | Reg-VQ 1 | Reg-VQ 2 | Reg-VQ 3 |
|-------|----------|----------|----------|



Figure 10. Semantic image synthesis on ADE20K with regularized quantizer and auto-regressive model.

Figure 11. Semantic image synthesis on CelebA-HQ with regularized quantizer and auto-regressive model.

| Input | Reg-VQ 1 | Reg-VQ 2 | Reg-VQ 3 |
|-------|----------|----------|----------|

A green bird with yellow belly on the tree.

A black bird with white breast and black crown.

A grey bird has white belly and grey wings.

A blue bird with yellow throat on the tree.

A yellow bird with yellow belly and grey wings.

Figure 12. Text-to-image synthesis on CUB-200 with regularized quantizer and diffusion model.

| Input | Reg-VQ 1 | Reg-VQ 2 | Reg-VQ 3 |
|---|---|---|---|

A red bus is running on street in the day.

There is a group of giraffes in the wild.

A man is riding a skateboard on the street.

Someone walking with a surfboard on some sand.
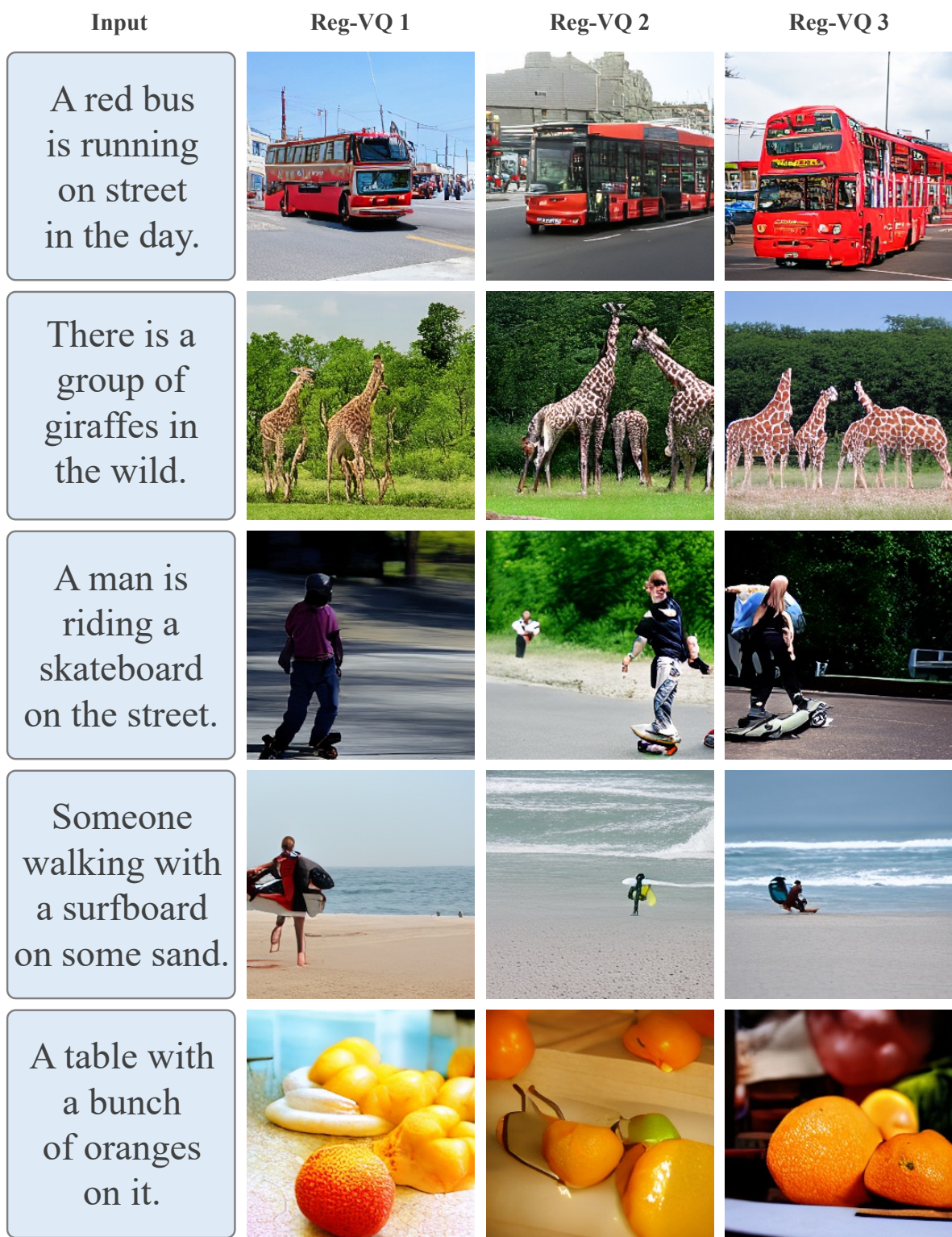
A table with a bunch of oranges on it.



Figure 13. Text-to-image synthesis on MS-COCO with regularized quantizer and diffusion model.