# RAQ-VAE: RATE-ADAPTIVE VECTOR-QUANTIZED VARIATIONAL AUTOENCODER

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Vector Quantized Variational AutoEncoder (VQ-VAE) is an established technique in machine learning for learning discrete representations across various modalities. However, its scalability and applicability are limited by the need to retrain the model to adjust the codebook for different rate requirements or encoding efficiency. We introduce the Rate-Adaptive VQ-VAE (**RAQ-VAE**) framework, which addresses this challenge with two novel discrete (codebook) representation methods: a model-based approach using a clustering technique for existing pre-trained VQ-VAE models, and a data-driven approach utilizing a sequence-to-sequence (Seq2Seq) model for variable-rate codebook generation. Our experiments demonstrate that RAQ-VAE achieves effective reconstruction performance across multiple rates, often outperforming conventional fixed-rate VQ-VAE models. This work enhances the adaptability and performance of VQ-VAEs, with broad applications in data reconstruction, generation, and computer vision tasks.

023

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

025 026

Vector quantization (VQ) (Gray, 1984) is a fundamental technique for learning discrete representations for various tasks (Krishnamurthy et al., 1990; Gong et al., 2014; Van Niekerk et al., 2020)
in the field of machine learning. The Vector Quantized Variational AutoEncoder (VQ-VAE) (Van Den Oord et al., 2017; Razavi et al., 2019), which extends the encoder-decoder structure of the Variational Autoencoder (VAE) (Kingma & Welling, 2013; Rezende & Viola, 2018), introduces discrete latent representations that have proven effective across multiple domains, including computer vision (Razavi et al., 2019; Esser et al., 2021), audio (Dhariwal et al., 2020; Yang et al., 2023; Tseng et al., 2023), and speech (Kumar et al., 2019; Xing et al., 2023). These successes are attributed to the inherently discrete nature of the data in these domains, which makes VQ well suited to learning complex inference and prediction tasks.

Recent developments have further enhanced VQ-based discrete representation learning by integrat-037 ing it with deep generative models, such as Generative Adversarial Networks (GANs) (Esser et al., 2021) and Denoising Diffusion Probabilistic Models (DDPMs) (Cohen et al., 2022; Gu et al., 2022; Yang et al., 2023). As VQ-VAE is integrated into these diverse generative frameworks, its utility and 040 applicability in various tasks are becoming increasingly evident. However, even with this success, 041 the scalability of the codebook-driven quantization process poses a significant challenge, further mo-042 tivating our approach. With the proliferation of large datasets and the demand for real-time process-043 ing, VQ-based architectures struggle with the computational complexity associated with dynamic 044 compression, including the need to retrain models to adjust computational loads. Consequently, addressing the scalability of the VQ process is crucial to fully realizing the potential of VQ-VAE, especially in integrating it with large-scale generative VQ models (Yu et al., 2022). 046

To address the issues, Li et al. (2023) proposed a method to resize the codebook without retraining the publicly available VQ models by applying hyperbolic embeddings to enhance the codebook vector with co-occurrence information and reordering the improved codebook with a Hilbert curve. Another approach to achieve more comprehensive codebook representation, the use of multi-codebook has been an ongoing challenge to achieve richer representations for different tasks (Guo et al., 2022).
Malka et al. (2023) designed and learned a nested codebook based on progressive learning to support different quantization levels. Guo et al. (2023) proposed a framework for predicting codebook indexes generated from embeddings of student models using multi-codebook vector quantization

by reformulating teacher label generation as a codec problem in knowledge distillation. Recently, Huijben et al. (2024) focused on unsupervised codebook generation based on residual quantization by studying the vector quantizer itself. However, addressing these issues through multi-codebook or residual quantization generally entails substantial changes to the existing well-established structure of VQ-VAEs, or face a reduction in the resolution of the quantized feature map.

To this end, we propose a Rate-Adaptive VQ-VAE (RAQ-VAE) framework that allows discrete 060 representation at various rates with a single VQ-VAE model. First, we propose model-based RAQ-061 VAE that can use the existing VQ-VAE model to obtain rate-adaptive VQ through a differential 062 k-means clustering (DKM) (Cho et al., 2021) algorithm and its inverse functionalization without 063 any additional parameters and retraining. Next, we present data-driven RAQ-VAE with Sequenceto-Sequence (Seq2Seq) (Sutskever et al., 2014) model for rate-adaptive codebook generation. The 064 data-driven RAQ-VAE can achieve discrete representation at any desired rate through the Seq2Seq 065 model and approaches or partially outperforms the separately trained conventional VQ-VAE model. 066 Our framework addresses the challenge of needing multiple VQ-VAE models for different compres-067 sion rates, especially in large-scale computer vision tasks that require high-capacity representations. 068 Additionally, it can be seamlessly integrated into various VQ applications without requiring signifi-069 cant modifications to the existing VQ-VAE structure.

- 071 Our contributions are summarized as follows:
  - We introduce the RAQ-VAE framework with two VQ codebook representation methods: *model-based* RAQ-VAE, utilizing an existing trained VQ-VAE model, and *data-driven* RAQ-VAE, combining Seq2Seq model with VQ-VAE architecture.
  - We propose *model-based* RAQ-VAE, which adapts the codebook of a VQ-VAE model using a dynamic codebook clustering method, allowing the quantizer to adjust the rate without retraining.
  - We propose *data-driven* RAQ-VAE that generates a rate-adaptive codebook via a Seq2Seq model. This approach uses a single codebook and a training method, *cross-forcing*, to train recurrent networks to generate codebooks at different rates.
    - Our experiments demonstrate that a single RAQ-VAE model achieves or even outperforms the performance of multiple VQ-VAE models trained at fixed rates, using the same encoderdecoder architecture.

#### 2 BACKGROUND

072

073

074

075

076

077

078

079

080

081

082

084 085

097 098

103 104 105

Vector-Quantized Autoencoder VQ-VAEs (Van Den Oord et al., 2017; Razavi et al., 2019) can successfully represent meaningful features that span multiple dimensions of data space by discretizing continuous latent variables to the nearest code vector in its codebook. In VQ-VAE, learning of discrete representations is achieved by quantizing the encoded latent variables to their nearest neighbors in a trainable codebook and decoding the input data from the discrete latent variables. To represent the data  $\mathbf{x} \in \mathbb{R}^{3 \times H \times W}$  from dataset  $\mathcal{D}$  discretely, a codebook e consisting of K learnable code vectors  $\{e_i\}_{i=1}^K \subset \mathbb{R}^d$  is employed. The quantized discrete latent variable  $\mathbf{z}_q(\mathbf{x}|\mathbf{e})$  is decoded to reconstruct the data  $\mathbf{x}$ . The quantizer Q modeled as deterministic encoder  $f_{\phi}$  to  $\mathbf{z}_q(\mathbf{x}|\mathbf{e})$  by finding the nearest neighbor in the D-dimensional codebook  $\mathbf{e} = \{e_i\}_{i=1}^K$  as

$$\mathbf{z}_{q}(\mathbf{x}|\mathbf{e}) = Q\left(f_{\phi}(\mathbf{x})|\mathbf{e}\right) = \operatorname*{arg\,min}_{e_{i} \in \{e_{i}\}_{i=1}^{K}} \left\|f_{\phi}(\mathbf{x}) - e_{i}\right\|.$$
(1)

The quantized representation is fixed to  $\log_2 K$  bits for the index *i* of the selected code vector  $e_i$ of the codebook of size *K*. The deterministic decoder  $f_{\theta}$  reconstruct the data **x** from the quantized discrete latent variable  $\mathbf{z}_q(\mathbf{x}|\mathbf{e})$  as  $\hat{\mathbf{x}} = f_{\theta}(\mathbf{z}_q(\mathbf{x}|\mathbf{e})|\mathbf{e}))$ . During the training process, the encoder  $f_{\phi}$ , decoder  $f_{\theta}$ , and codebook **e** are jointly optimized to minimize the loss  $\mathcal{L}_{VO}(\phi, \theta, \mathbf{e}; \mathbf{x}) =$ 

$$\underbrace{\log p_{\theta}(\mathbf{x}|\mathbf{z}_{q}(\mathbf{x}|\mathbf{e}))}_{\mathcal{L}_{\text{recon}}} + \underbrace{\left|\left|sg\left[f_{\phi}(\mathbf{x})\right] - \mathbf{z}_{q}(\mathbf{x}|\mathbf{e})\right|\right|_{2}^{2}}_{\mathcal{L}_{\text{embed}}} + \underbrace{\beta\left|\left|sg\left[\mathbf{z}_{q}(\mathbf{x}|\mathbf{e})\right] - f_{\phi}(\mathbf{x})\right|\right|_{2}^{2}}_{\mathcal{L}_{\text{commit}}}$$
(2)

106 where sg[·] is the *stop-gradient* operator. The  $\mathcal{L}_{recon}$  is the reconstruction loss between the input 107 data x and the reconstructed decoder output  $\hat{x}$ . The two  $\mathcal{L}_{embed}$  and  $\mathcal{L}_{commit}$  apply only to codebook variables and encoder weight with a weighting hyperparameter  $\beta$  to prevent fluctuations from



Figure 1: Model-based RAQ-VAE (*left*): the model-based approach clusters the codebook e of a trained VQ-VAE model with separate tasks for reducing or increasing to the adapted codebook  $\tilde{e}$  and applies it to the model. Data-driven RAQ-VAE (*right*): the data-driven approach trains a Seq2Seqbased codebook adaptation procedure utilizing the baseline VQ-VAE model with data, where the gradient flow of the codebook passes through the Seq2Seq model.

one code vector to another. Since the quantization process is non-differentiable, the codebook loss is typically approximated via a straight-through gradient estimator (Bengio et al., 2013), such as  $\partial \mathcal{L}/\partial f_{\phi}(\mathbf{x}) \approx \partial \mathcal{L}/\partial \mathbf{z}_q(\mathbf{x})$ . Both conventional VAE (Kingma & Welling, 2013) and VQ-VAE (Van Den Oord et al., 2017) have objective functions consisting of the sum of reconstruction error and latent regularization. To improve performance and convergence rate, exponential moving average (EMA) update is usually applied for the codebook optimization (Van Den Oord et al., 2017; Razavi et al., 2019) (for more details in supplementary material A.1).

138 Seq2Seq The sequence-to-sequence (Seq2Seq) (Sutskever et al., 2014) model is widely used in 139 sequence prediction tasks such as language modeling and machine translation (Dai & Le, 2015; Lu-140 ong et al., 2016; Ranzato et al., 2016). The model employs an initial LSTM, called the encoder, 141 to process the input sequence sequentially and produce a substantial fixed-dimensional vector rep-142 resentation, called the context vector. The output sequence is then derived by a further LSTM, the decoder. In particular, the decoder is conditioned on the input sequence, distinguishing it as a dis-143 tinct component within the architecture. During training, the Seq2Seq model typically uses *teacher* 144 forcing (Williams & Zipser, 1989), where the target sequence is provided to the decoder at each time 145 step, instead of the decoder using its own previous output as input. This method helps the model 146 converge faster by providing the correct context during training. 147

148 149

150

130

#### 3 RATE-ADAPTIVE VQ-VAE

151 Although VQ-VAE has been successfully applied to various domains, it still faces scalability limi-152 tation. In particular, the common fixed-rate VO-VAE model requires modifying the codebook size 153 K when processing different datasets (see (Razavi et al., 2019; Esser et al., 2021), using as many as 16384 and as few as 512 codebook sizes are used). Furthermore, adjusting the computational load 154 requires retraining the model, which poses additional challenges. To overcome these limitations, we 155 introduce two Rate-Adaptive VQ-VAE (RAQ-VAE) frameworks, which can adjust the rate of VQ-156 VAE through increasing or decreasing of the codebook size K. The outline of the RAQ-VAE frame-157 work is shown in Figure 1. RAQ-VAE builds upon a codebook mapping  $\Psi : (\mathbb{R}^d)^{\times K} \longrightarrow (\mathbb{R}^d)^{\times K}$ 158 for any integer  $\widetilde{K} \in \mathbb{N}$  that can be either lower, i.e.,  $\widetilde{K} < K$ , or higher, i.e.,  $\widetilde{K} > K$ , than the 159 original codebook size K. We design the mapping in two ways: (i) model-based RAQ-VAE; (ii) 160 data-driven RAQ-VAE. Model-based RAQ-VAE (Sec. 3.1) can obtain rate-adaptive VQ through 161 differentiable k-means clustering (DKM) (Cho et al., 2021) algorithm without any additional pa162 rameters. In addition, data-driven RAQ-VAE (Sec. 3.2) is an offline-trained RAQ-VAE method that 163 adopts the codebook generative Sequence-to-Sequence (Seq2Seq) (Sutskever et al., 2014) model. 164

3.1 MODEL-BASED RATE-ADAPTIVE VQ-VAE

167 Previous attempts (Łańcucki et al., 2020; Tjandra et al., 2019; Zheng & Vedaldi, 2023) have pro-168 posed enhancing codebook learning by periodically clustering the codebook during model training. In contrast, we propose a model-based rate-adaptive VQ-VAE that performs online codebook clus-169 170 tering after the model has been trained. By loading a VQ-VAE model trained with the original codebook  $\tilde{\mathbf{e}}$  and dynamically clustering the codebook to the adapted codebook  $\tilde{\mathbf{e}}$ . This allows the 171 vector quantizer to adapt to nuanced patterns within the overall model, providing flexibility and 172 scalability (See Figure 1).

173 174

185

186

187

188 189 190

192

193

194

196 197

199 200

201

202 203

204

205 206 207

166

**Codebook Clustering** To achieve the desired rate for the adapted codebook size  $\tilde{K}$  (=  $|\tilde{\mathbf{e}}|$ ), we 175 derive the clustered codebook e from the original codebook e. Details of the codebook clustering 176 formulation are provided in supplementary material A.2. To ensure that the clustering process is 177 effectively integrated into the trained VQ-VAE model, we employ a differentiable k-means clus-178 tering (DKM) algorithm (Cho et al., 2021). This algorithm, originally proposed for DNN model 179 compression, uses an attention-based weight clustering method. We use the DKM algorithm for VQ codebook clustering, focusing on the fine-tuning of clustered codebooks and VQ-VAE model 181 architectures. Additionally, we utilize DKM for codebook incrementation (inverse functionalization 182 process) to handle scenarios requiring an increase in codebook size. 183

**Reducing the Rate**  $(\tilde{K} < K)$  In the rate reduction task, DKM can perform iterative differentiable codebook clustering on  $\widetilde{K}$  clusters. Let C represent the cluster centers and vector e represent the original codebook. The DKM algorithm for VQ codebook operates as follows:

- Initialize a centroid C = {c<sub>j</sub>}<sup>˜</sup><sub>j=1</sub> either by randomly selected ˜K codebook vectors from e or using k-means++. The last known C from the previous batch is used for all following iterations.
  - Calculate the distance between the original codebook vector  $e_i$  and initialized centroid  $c_i$ using Euclidean distance as the distance metric  $D_{i,j} = -f(e_i, c_j)$  with its matrix **D**.
  - To obtain the attention matrix A, derive each row of A where  $A_{i,j} = \frac{\exp\left(\frac{D_{i,j}}{\tau}\right)}{\sum_k \exp\left(\frac{D_{i,k}}{\tau}\right)}$  represents the attention from

sents the attention from  $e_i$  and  $c_j$  with a softmax temperature  $\tau$ .

- Get a centroid candidate  $\widetilde{\mathbf{C}} = \{\widetilde{c}_j\}_{j=1}^{\widetilde{K}}$  by summing all the attentions for each centroid by computing  $\widetilde{c}_j = \frac{\sum_i A_{i,j} e_i}{\sum_i A_{i,j}}$  and update  $\mathbf{C}$  with  $\widetilde{\mathbf{C}}$ .
- Repeat this process until  $|\mathbf{C} \mathbf{C}| \le \epsilon$  at which point DKM has converged or the iteration limit reached, then compute AC to get  $\tilde{e}$ .

The above iterative process can be summarized as follows:

$$\tilde{\mathbf{e}} = \underset{\tilde{\mathbf{e}}}{\operatorname{arg\,min}} \mathcal{L}_{\text{DKM}}(\mathbf{e}; \tilde{\mathbf{e}}) = \underset{\mathbf{C}}{\operatorname{arg\,min}} \left| \mathbf{C} - \mathbf{A}\mathbf{C} \right| = \underset{\mathbf{C}}{\operatorname{arg\,min}} \sum_{j=1}^{K} \left| c_j - \frac{\sum_i A_{i,j} e_i}{\sum_i A_{i,j}} \right|$$
(3)

~.

208 In (Cho et al., 2021), the authors implemented DKM for soft-weighted cluster assignment and hard-209 ness can be enforced to provide weighted clustering constraints. In the softmax operation, the tem-210 perature  $\tau$  can be used to control the level of hardness. At the end of the DKM process, we use the 211 last attention matrix A to snap each codebook vector to the nearest centroid and finish clustering the codebook. 212

213

**Increasing the Rate** ( $\tilde{K} > K$ ) While k-means clustering is effective for compressing code vec-214 tors, it has algorithmic limitations that prevent the augmentation of additional codebooks. To ad-215 dress this, we introduce the inverse functional DKM (IKM), a technique for increasing the number of codebooks. This iterative method aims to approximate the posterior distribution of an existing gener ated codebook. We use maximum mean discrepancy (MMD) to compare the distribution difference
 between the base codebook and the clustered generated codebook, where MMD is a kernel-based
 statistical test technique that measures the similarity between two distributions (Gretton et al., 2012).

Assuming the original codebook vector  $\mathbf{e}$  of size K already trained in the baseline VQ-VAE, the process of generating the codebook  $\tilde{\mathbf{e}}$  using the IKM algorithm is performed as follows:

- 223 224
- Initialize a *d*-dimensional adapted codebook vector  $\tilde{\mathbf{e}} = \{\tilde{e}_i\}_{i=1}^{\tilde{K}}$  as  $\tilde{\mathbf{e}} \sim \mathcal{N}(0, d^{-\frac{1}{2}} \boldsymbol{I}_{\tilde{K}})$
- Cluster  $\tilde{\mathbf{e}}$  via the DKM process (equation 3):  $g_{\text{DKM}}(\tilde{\mathbf{e}}) = \arg\min_{q_{\text{DKM}}(\tilde{\mathbf{e}})} \mathcal{L}_{\text{DKM}}(\tilde{\mathbf{e}}; g_{\text{DKM}}(\tilde{\mathbf{e}})).$
- 225 226 227

228

248

- Calculate the MMD between the true original codebook e and the DKM clustered  $g_{\text{DKM}}(\tilde{\mathbf{e}})$ .
- Optimize  $\tilde{\mathbf{e}}$  to minimize the MMD objective  $\mathcal{L}_{\text{IKM}}(\mathbf{e}; \tilde{\mathbf{e}}) = \text{MMD}(\mathbf{e}, g_{\text{DKM}}(\tilde{\mathbf{e}})) + \lambda ||\tilde{\mathbf{e}}||^2$ .

229 where  $\lambda$  is the regularization parameter controlling the strength of the L2 regularization term. The 230 IKM process can be summarized as  $\tilde{\mathbf{e}} = \arg \min \mathcal{L}_{IKM}(\mathbf{e}; \tilde{\mathbf{e}})$ . Since DKM does not block gradient 231 flow, we easily can update the codebook  $\tilde{\mathbf{e}}$  using stochastic gradient descent (SGD) as  $\tilde{\mathbf{e}} = \tilde{\mathbf{e}} -$ 232  $\eta \nabla \mathcal{L}_{IKM}(\mathbf{e}, \tilde{\mathbf{e}})$ . With DKM and IKM, the generated codebook  $\tilde{\mathbf{e}}$  can be used to quantize the encoded 233 vector as  $\mathbf{z}_q(\mathbf{x}|\tilde{\mathbf{e}})$  at different rates without adding any model parameters to the trained VQ-VAE. 234 Since DKM does not block gradient flow, it is easy to change the codebook cluster assignments 235 even during offline and online training. During offline training, the clusters that are best suited 236 in terms of VQ task loss are adopted. Although we do not focus on using multi-codebook with 237 DKM (aim to leverage rate-adaptive codebook after trained), a multi-codebook VQ-VAE model can 238 be easily implemented by tuning K with DKM during training and hierarchically optimizing the 239 multi-codebook clusters with the model. 240

241<br/>2423.2DATA-DRIVEN RATE-ADAPTIVE VQ-VAE

Seq2Seq models (Sutskever et al., 2014) have been widely used in machine translation to handle
variable output sequences, where the length of sentences can differ significantly between languages.
Inspired by this, we propose a Seq2Seq-based approach to generate rate-adaptive codebooks within
the VQ-VAE framework. This section introduces the data-driven RAQ-VAE, which integrates a
learning vector quantization layer with Seq2Seq model.

**Overview** As shown in Figure 1, data-driven RAQ-VAE is constructed with a deterministic encoder-decoder pair, a trainable original codebook e, and Seq2Seq model. The adapted codebook  $\tilde{e}$ is generated by the Seq2Seq model from the original codebook e. Data-driven RAQ-VAE hierarchically quantizes the continuous latent representation  $f_{\phi}(x)$  of data x into  $\mathbf{z}_q(\mathbf{x}|\mathbf{\hat{e}})$  and  $\mathbf{z}_q(\mathbf{x}|\mathbf{\hat{e}})$  via e and  $\tilde{e}$ , respectively. Building on the conventional VQ-VAE architecture, the data-driven RAQ-VAE learns the encoder-decoder pair while training the codebook e and its generative process  $G_{\psi}$ .

255 **Codebook Encoding** The rate-adaptive codebook generation procedure,  $G_{\psi}$ , leverages LSTM 256 cells in the Seq2Seq model to dynamically generate an adapted codebook e from the original code-257 book e. The first step is to initialize the target codebook size K. During training, the data-driven 258 RAQ-VAE is trained with arbitrary codebook sizes K. In the test phase, the Seq2Seq model gen-259 erates the adapted codebook  $\tilde{\mathbf{e}}$  at the desired rate specified by the user. This initialization sets the 260 foundation for the encoding and decoding steps in Algorithm 1. Each vector of the original code-261 book  $e_i$  is sequentially encoded by a set of LSTM cells. The hidden and cell states (h, c) capture 262 the contextual information of each base codebook vector. 263

Codebook Decoding via cross-forcing The goal of Seq2Seq codebook generation is to reflect as
 much information as possible from the original codebook while generating a usable codebook for the
 VQ-VAE decoder. However, existing Seq2Seq training methods, such as teacher forcing (Williams
 & Zipser, 1989), may not be suitable when the target adapted codebook ẽ consists of sequences that
 are much longer than the original codebook. Therefore, we propose cross-forcing, a hybrid approach
 combining teacher forcing and free running in professor forcing (Lamb et al., 2016). This is feasible
 because, unlike typical sequence prediction tasks, the order of the codebooks does not significantly

| <b>Input:</b> Original codebook $\mathbf{e} = \{e_i\}_{i=1}^K$  | <b>Input:</b> x (batch of training data)  |
|---|---|
| <b>Output:</b> Adapted codebook $\tilde{\mathbf{e}} = {\tilde{e}_i}_{i=1}^{\tilde{K}}$  | for $\mathbf{x} \in$ train dataset $\mathcal{D}$ do   |
| Initialize adapted codebook size $\widetilde{K}$ ,<br>hidden $\mathbf{h} = \{h_i\}_{i=1}^K$ , and cell $\mathbf{c} = \{c_i\}_{i=1}^K$ . | $\mathbf{z}_q(\mathbf{x} \mathbf{e}) \leftarrow Q\left(f_{\phi}(\mathbf{x}) \mathbf{e}\right)$  |
| $\triangleright \text{ Codebook encoding}$  | $\triangleright$ Generate $\tilde{\mathbf{e}}$ from Seq2Seq model $G_{abc}$   |
| for $i = 1$ to K do   | $\mathbf{\tilde{e}} \leftarrow G_{\psi}(\mathbf{e})$ by Algorithm 1   |
| $h_i, c_i \leftarrow LSTM_\psi(e_i)$ end for  | $\triangleright$ Ouantize encoder output $f_{\phi}(\mathbf{x})$ with $\tilde{\mathbf{e}}$ .   |
| $\triangleright$ Codebook decoding via <i>cross-forcing</i>   | $\mathbf{z}_{q}(\mathbf{x} \tilde{\mathbf{e}}) \leftarrow Q\left(f_{\phi}(\mathbf{x}) \tilde{\mathbf{e}} ight)$   |
| for $i = 1$ to K do   |   |
| If $i < 2K$ and <i>i</i> is odd then<br>$\tilde{e}_i \leftarrow LSTM_{ib}(e_i, h, c)$   | $\mathbf{\hat{x}_{e}}, \mathbf{\hat{x}_{\tilde{e}}} \leftarrow f_{\theta}(\mathbf{z}_{q}(\mathbf{x} \mathbf{e})), f_{\theta}(\mathbf{z}_{q}(\mathbf{x} \mathbf{e}))$<br>Compute $\mathcal{L}_{VO}$ by equation 2. |
| else $\psi(z_1, \dots, z_n)$  | Compute $\mathcal{L}_{RAO}$ by equation 4.  |
| $\tilde{e}_i \leftarrow LSTM_{\psi}(\tilde{e}_i, \boldsymbol{h}, \boldsymbol{c})$   | $\phi, \theta, \mathbf{e} \leftarrow \text{Update}(\mathcal{L}_{VO})$   |
| end if  | $\phi, \theta, \psi, \mathbf{e} \leftarrow \text{Update}(\hat{\mathcal{L}}_{RAQ})$  |
| end for   | end for   |
| <b>Return:</b> $\tilde{\mathbf{e}} = G_{\psi}(\mathbf{e})$  | <b>Return:</b> $f_{\phi}, f_{\theta}, G_{\psi}, \mathbf{e}$   |

affect the outcome. In the decoding phase (as shown in Algorithm 1), teacher forcing is applied for odd steps that are less than twice the original codebook size  $(2\tilde{K})$ , using the base code vector  $(e_i)$  as input. For even steps and beyond, *free running* (using the previous time step decoder output as input) is performed to dynamically train the VQ-VAE decoder with the generated codebook. The codebook decoding via cross-forcing is a key component of the data-driven approach. This technique helps ensure stable codebook generation at different rates. We have empirically validated its effectiveness in Appendix A.4.5.

Training Procedure To train the data-driven RAQ-VAE, we jointly optimize the base VQ-VAE and RAQ-VAE objectives to learn a good representation of the original codebook e and the rateadaptive codebook generative process  $G_{\psi}$ . We formulate the constrained optimization  $\mathcal{L}_{RAQ}$  to jointly update  $G_{\psi}$  with  $f_{\phi}$ ,  $f_{\theta}$ , and e as  $\mathcal{L}_{VQ}(\phi, \theta, \mathbf{e}; \mathbf{x}) \geq \mathcal{L}_{RAQ}(\phi, \theta, \psi, \mathbf{e}; \mathbf{x}) =$ 

$$\log p_{\theta} \left( \mathbf{x} | \mathbf{z}_q(\mathbf{x} | G_{\psi}(\mathbf{e})) \right) + \left| \left| \sup \left[ f_{\phi}(\mathbf{x}) \right] - \mathbf{z}_q(\mathbf{x} | G_{\psi}(\mathbf{e})) \right| \right|_2^2 + \beta \left| \left| \sup \left[ \mathbf{z}_q(\mathbf{x} | G_{\psi}(\mathbf{e})) \right] - f_{\phi}(\mathbf{x}) \right| \right|_2^2.$$
(4)

where sg[·] is the *stop-gradient* operator. The data-driven RAQ-VAE jointly minimizes  $\mathcal{L}_{VQ}$  (equation 2) and  $\mathcal{L}_{RAQ}$  (equation 4). Back-propagating  $\mathcal{L}_{VQ}$  induces the same gradient flows as the base VQ-VAE. Additionally, back-propagating  $\mathcal{L}_{RAQ}$  induces a gradient flow to the Seq2Seq model, resulting in effective codebook generation. The overall training procedure for the proposed data-driven RAQ-VAE is summarized in Algorithm 2. During training, the Seq2Seq model dynamically generates codebooks and adapts to different rates at each training iteration.

309 310 311

303 304

305

306

307

308

290

#### 4 RELATED WORK

312 313

VQ-VAE and its Improvements The VQ-VAE (Van Den Oord et al., 2017) has inspired numerous 314 developments since its inception. Łańcucki et al. (2020); Williams et al. (2020); Zheng & Vedaldi 315 (2023) proposed codeword reset and online clustering methods to address the problem of *codebook* 316 collapse (Takida et al., 2022), thereby increasing the training efficiency of the codebook. Tjandra 317 et al. (2019) introduced a conditional VQ-VAE that generates magnitude spectrograms for target 318 speech using a multi-scale codebook-to-spectrogram inverter given the VQ-VAE codebook. SQ-319 VAE (Takida et al., 2022) incorporated stochastic quantization and a trainable posterior categorical 320 distribution to enhance VQ-VAE performance, while Vuong et al. (2023) proposed VQ-WAE, based 321 on SQ-VAE, using Wasserstein distance to ensure a uniform distribution of discrete representations. Several works have introduced substantial structural changes to VQ-VAE. Lee et al. (2022) proposed 322 a two-step framework with Residual Quantized (RQ) VAE and RQ-Transform to generate high-323 resolution images using a single shared codebook. Mentzer et al. (2023) replaced VQ with Finite

Scalar Quantization (FSQ) to tackle codebook collapse. However, unlike previous works, we focus
 on achieving rate-adaptive VQ-VAE within a largely unchanged quantization scheme and VQ-VAE
 model architecture to improve its scalability for application not only to basic VQ-VAE models but
 also to its advanced models.

- Variable-Rate Neural Image Compression Several studies have proposed variable-rate learn-330 ing image compression frameworks based on different neural network architectures. Yang et al. 331 (2020); Choi et al. (2019); Cui et al. (2020) introduced frameworks based on autoencoders, con-332 ditional autoencoders, and VAE structures, respectively. Variable-rate image compression has also been achieved in studies such as Song et al. (2021), which uses models based on the Spatial Fea-333 ture Transform (SFT) for compression, and Johnston et al. (2018), which employs recurrent neural 334 networks (RNNs) to achieve variable-rate compression by evaluating the distortion of individual 335 patches to compute a weighted distortion. Duong et al. (2023) proposed learned transforms and 336 entropy coding to enhance the linear transforms in existing codecs by systematizing the process into 337 a single model that follows the rate-distortion curve. However, the integration of variable-rate im-338 age compression within the VQ-VAE framework remains an open question. Unlike these studies, 339 our work focuses on embedding variable-rate compression directly into the VQ-VAE framework, 340 maintaining the benefits of VQ while enhancing scalability and adaptability.
  - 5 EXPERIMENTS
- 342 343 344

341

DEAPERIMENTS

345ImplementationTo demonstrate the advancement of the proposed RAQ-VAE, we adapt the con-<br/>ventional VQ-VAE (Van Den Oord et al., 2017) and the two-level hierarchical VQ-VAE (VQ-VAE-<br/>2) (Razavi et al., 2019) as baselines. We perform empirical evaluations on vision datasets: CIFAR10<br/> $(32 \times 32)$  (Krizhevsky et al.) and CelebA ( $64 \times 64, 128 \times 128$ ) (Liu et al., 2015) for quantitative<br/>evaluation, and ImageNet ( $256 \times 256$ ) (Russakovsky et al., 2015) for qualitative evaluation. We<br/>designed RAQ-VAE to adapt the conventional VQ-VAE and its improved model structures (Tjandra<br/>et al., 2019; Ott et al., 2019; Esser et al., 2021; Ramesh et al., 2021) to achieve multiple rates within<br/>a single model.

Architecture We use identical architecture and parameters for all methods, setting the default codeword (discrete latent) embedding dimension *d* to 64 for CIFAR10 and CelebA, and to 128 for ImageNet. The codebook sizes range from 16 to 1024 for CIFAR10, 32 to 2048 for CelebA, and 128 to 4096 for ImageNet, with conventional VQ-VAE models trained on 'power of 2' sizes and RAQ-VAE models set to the middle of the range for both model-based and data-driven approaches. Details of the experimental settings are provided in supplementary material A.3.

359

352

**Evaluation Metrics** We quantitatively evaluated our method using peak-signal-to-noise-ratio 360 (PSNR), structural similarity index measure (SSIM), reconstructed Fréchet inception distance (rFID) 361 (Heusel et al., 2017), and codebook perplexity. PSNR measures the ratio between the maximum 362 possible power of a signal and the power of the corrupted noise affecting data fidelity (Korhonen 363 & You, 2012). SSIM assesses structural similarity between two images (Wang et al., 2004; Brunet 364 et al., 2011). rFID assesses the quality of reconstructed images by comparing the distribution of 365 features extracted from the test data with that of the original data. Codebook perplexity, defined as 366  $e^{-\sum_{i}^{\widetilde{K}} p_{e_i} \log p_{e_i}}$  where  $p_{e_i} = \frac{N_{e_i}}{\sum_{i}^{\widetilde{K}} N_{e_i}}$  and  $N_{e_i}$  represents the encoded number for latent representa-367 tion with codebook  $e_i$ , indicates a uniform prior distribution when the perplexity value reaches the 368 369 codebook size (K), meaning all codebooks are used equally.

370 371

372

#### 5.1 MAIN RESULTS ON VISION TASKS

**Quantitative Evaluation** We empirically compare our RAQ-VAE models with the conventional VQ-VAE (Van Den Oord et al., 2017) for image reconstruction performance. We trained and evaluated each VQ-VAE with different codebook sizes (K) as a quantitative baseline, and then validated RAQ-VAE by adapting the rate (by adjusting  $\tilde{K}$ ) on a single model-based and data-driven RAQ-VAE model. Figure 2 shows the results, evaluated on the CIFAR10 and CelebA ( $64 \times 64$ ) datasets. Under same compression rate and network architecture, all proposed RAQ-VAE models achieve



Figure 2: **Reconstruction performance** at different rates (adapted codebook sizes) evaluated on (a) CIFAR-10 and (b) CelebA. Higher values are better for PSNR, SSIM, and codebook perplexity, while lower values are better for rFID. The black lines represent separate VQ-VAE models trained individually for each codebook size  $\tilde{K}$ . The colored lines represent a single RAQ-VAE model initially trained at a original codebook size K and later adapted to different codebook sizes  $\tilde{K}$ . The shaded area indicates the 95.45% confidence interval based on 4 runs with different training seeds.

performance close to that of multiple VQ-VAE models. When increasing the rate (codebook size), 404 the data-driven RAQ-VAE achieves slightly lower results for the PSNR and SSIM metrics but bet-405 ter results in terms of rFID score, which evaluates perceptual image quality at the dataset level. In 406 particular, the perplexity of the conventional VQ-VAE models shows low scores on CelebA, but the 407 proposed data-driven RAQ-VAE performs better in terms of perplexity and rFID, especially at high 408 bits per pixel (bpp). The model-based RAQ-VAE performs poorly overall, but in the task of reducing 409 the rate, it achieves intermittently more reliable results on CIFAR10. Our proposed method is highly 410 portable and reduces model complexity, considering the resources invested in each single VQ-VAE, 411 since RAQ-VAE covers multiple fixed-rate VQ-VAE models with only a single model. (The model 412 complexity of the baseline VQ-VAEs and our RAQ-VAEs are provided in A.3.3.) 413

- 414 Qualitative Evaluation For qualitative evaluation, we compare a single data-driven RAQ-VAE 415 with VQ-VAEs trained at different rates (0.4375 bpp to 0.75 bpp) on ImageNet ( $256 \times 256$ ). As seen in Figure 3, the VO-VAEs (in the first row) are trained for each rate show that the quality 416 decreases as the rate decreases, which is consistent with the results observed in the quantitative eval-417 uation. When randomly selecting codebooks from a VQ-VAE model trained with K = 4096 (in 418 the second row), we observe significant color changes, particularly at 0.5 bpp, where reconstruc-419 tions retain structural similarity but show color distortions. However, the data-driven RAQ-VAE 420 (K = 512), trained on a low-rate base codebook (0.5625 bpp), preserves the high-level semantic 421 features and colors of the input image well with only a single model trained on the low-rate base 422 codebook. Notably, it recovers fine details, like the cat's whiskers, far better than reconstructions 423 using randomly selected codebooks. Although the image quality declines slightly at the lowest bpp, 424 future work combining RAQ-VAE with advanced priors, such as PixelCNN (van den Oord et al., 425 2016) or PixelSNAIL (Chen et al., 2018), could further enhance the fidelity of generated images. 426 Additional reconstructions can be found in the supplementary material A.4.4.
- 427 428

429

403

5.2 DETAILED ANALYSIS

Codebook Usability Following the observations of previous works (Wu & Flierl, 2020; Takida et al., 2022; Vuong et al., 2023), we note that as the codebook size increases, the codebook perplexity of data-driven RAQ-VAE also increases, leading to better reconstruction performance. In most VQ-

471

480



| Method                                      | Dataset    | к               | Adapted Codebook Size $\widetilde{K}$ |       |                     |                     |        |        |
|---|------------|-----------------|---------------------------------------|-------|---------------------|---------------------|--------|--------|
| includu                                     | Dataset IX |                 | 2048                                  | 1024  | 512                 | 256                 | 128    | 64     |
| VQ-VAE-2 (Razavi et al., 2019)              | CelebA     | $\widetilde{K}$ | $10.11^{\S}$                          | 11.53 | 14.10               | $14.74^{\$\$}$      | 17.95  | 20.95  |
| VQ-VAE-2 <sup>†</sup> (Razavi et al., 2019) | CelebA     | 2048            | $10.11^{\$}$                          | 18.96 | 24.36               | 37.34               | 148.88 | 215.41 |
| Model-based RAQ* (Ours)                     | CelebA     | 256             | 13.76                                 | 13.68 | 14.12               | 14.74 <sup>§§</sup> | 20.29  | 53.46  |
| Data-driven RAQ* (Ours)                     | CelebA     | 256             | 13.48                                 | 13.65 | 14.18               | 15.88               | 19.33  | 25.24  |
|   |            |                 | 512                                   | 256   | 128                 | 64                  | 32     | -      |
| VQGAN (Esser et al., 2021)                  | ImageNet   | $\widetilde{K}$ | $13.71^{\$}$                          | 14.38 | $18.96^{\S\S}$      | 22.48               | -      | -      |
| VQGAN <sup>†</sup> (Esser et al., 2021)     | ImageNet   | 512             | $13.71^{\$}$                          | 17.26 | 23.94               | 35.92               | 60.79  | -      |
| Model-based RAQ** (Ours)                    | ImageNet   | 128             | 18.81                                 | 18.89 | 18.96 <sup>§§</sup> | 24.58               | 34.16  | -      |
| Data-driven RAQ** (Ours)                    | ImageNet   | 128             | 15.91                                 | 17.54 | 19.83               | 22.55               | 31.54  | -      |

Table 1: **Reconstruction performances** at different rates (according to  $\widetilde{K}$ ) on CelebA (128 × 128) test set and ImageNet (256 × 256) validation set. † uses a single model for reconstructions with randomly selected codebooks. \* denotes models trained with two-level hierarchical VQ-VAE (VQ-VAE-2) as in Razavi et al. (2019). \*\* denotes model trained with the **stage-1 VQGAN** as in Esser et al. (2021). § and §§ indicate results generated from the same model for the corresponding rates.

472 VAE frameworks, codebook perplexity is considered optimal when it approaches the codebook size, 473 effectively utilizing the available resources when the codebook size is limited. As demonstrated in 474 the main quantitative evaluation (see Figure 2), the data-driven RAQ-VAE outperforms conventional 475 VQ-VAE in terms of codebook perplexity at higher bits per pixel (bpp). This improvement highlights 476 the effectiveness of the Seq2Seq model in generating a codebook that the decoder can consistently 477 and efficiently utilize. The ability of data-driven RAQ-VAE to maintain high codebook perplexity 478 ensures better representation and reconstruction quality, proving its robustness in handling larger codebooks. 479

**Rate Adaptation** To demonstrate the rate adaptation performance, we validated RAQ-VAE by varying the adapted codebook size ( $\tilde{K}$ ). For the rate reduction task ( $\tilde{K} < K$ ), our experiments show that data-driven RAQ-VAE generally outperforms model-based RAQ-VAE in most aspects. However, on the CIFAR10, the model-based RAQ-VAE performs better at some rates. When a VQ-VAE model achieves high codebook perplexity, substantial performance can be achieved by simply clustering the codebook vectors (see more results in supplementary material A.4.1). For the rate

486 increasing task (K > K), a more challenging adaptation task, data-driven RAQ-VAE successfully 487 generated higher-rate codebooks, outperforming model-based RAQ-VAE and partially surpassing 488 conventional VQ-VAE models trained at the same codebook size. This capability was especially 489 pronounced on the CelebA dataset. For model-based RAQ-VAE, increasing the difference between 490 the original and adjusted codebook sizes resulted in noticeable performance degradation, exposing the limitations of the current implementation. However, the model-based approach can be advan-491 tageous for practitioners with limited computing resources as it allows them to just load and apply 492 codebook embeddings to the pre-trained VQ-VAE models without the need for additional training. 493 Although performance limitations remain, we suggest that adapting rates via the model-based ap-<u>191</u> proach could be another promising direction for future research. It is likely that large models such 495 as ViT-VQGAN (Yu et al., 2022) would experience even greater computational overhead compared 496 to CNN-based models, making this approach potentially beneficial. 497

497

**Applicability** To demonstrate the broader applicability of our methodology, we extend our ap-499 proach to the two-level hierarchical VQ-VAE (VQ-VAE-2) model (Razavi et al., 2019) and the 500 stage-1 VQGAN model (Esser et al., 2021) as the baseline models. The VQ-VAE-2 model is an ex-501 tension of the original VQ-VAE framework by incorporating a hierarchical structure that allows for 502 improved representation and reconstruction capabilities. The VQGAN enhances the encoding pro-503 cess in the first stage by incorporating adversarial and perceptual losses (Johnson et al., 2016; Zhang 504 et al., 2018), allowing for the generation of images with finer details. Table 1 provides the recon-505 struction performance according to the adapted codebook size K for the baseline models. Although 506 models trained at specific codebook sizes (first rows) achieve slightly better reconstruction, RAQ-VAE offers a flexible, efficient solution by covering multiple rates with a single adaptive model. In 507 most cases, the data-driven RAQ method outperforms the model-based approach. As shown in Fig-508 ure 3, training with a large codebook and then randomly selecting the codebook leads to significant 509 degradation when more than half of the codebook is removed. Applying our rate-adaptive quantiza-510 tion to VQ-VAE-2 and VQGAN not only preserves the performance of hierarchical or GAN-based 511 models but also provides the flexibility to adapt to different rates without retraining. This demon-512 strates that RAQ-VAE extends beyond VQ-VAE, offering a versatile solution for more advanced 513 VQ-based models, with significant potential in data reconstruction and generation tasks. 514

515

#### 6 CONCLUSION

516 517

We introduced the Rate-Adaptive VQ-VAE (RAQ-VAE) framework, which addresses the scalability 518 limitations of conventional VQ-VAEs through two novel codebook representation methods. Our ex-519 periments demonstrate that single RAQ-VAE model achieves superior reconstruction performance 520 across multiple rates without the need for retraining. The ability to dynamically adjust rates without 521 retraining makes it particularly beneficial for resource-constrained environments, simplifying model 522 deployment and management. This rate-adaptive capability provides significant flexibility for appli-523 cations that require dynamic compression levels, such as variable-rate image and video compression 524 (Xu et al., 2023) or real-time end-to-end communication systems (Park et al., 2020). Although 525 performance limitations remain, future work could further enhance stability and performance, in-526 creasing the overall value of our framework. With its proven versatility, RAQ-VAE has the potential 527 to drive significant advances in both the theoretical and practical fields of machine learning.

528 529

530

531

**Ethics Statement** RAQ-VAE is designed as a rate-adaptive extension of VQ-VAE and can be applied in all domains where VQ-based methods are used. As with all generative models, attention should be given to potential biases in the training data, as these can affect generated outputs. RAQ-VAE does not introduce any new ethical concerns beyond those inherent in VQ-VAE models.

532 533 534

**Reproducibility Statement** Appendix A.3 provides details of the experiments. The complete code necessary to reproduce our experiments is included in the supplementary material.

535 536

## 537 REFERENCES

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

- Dominique Brunet, Edward R Vrscay, and Zhou Wang. On the mathematical properties of the structural similarity index. *IEEE Transactions on Image Processing*, 21(4):1488–1499, 2011.
- XI Chen, Nikhil Mishra, Mostafa Rohaninejad, and Pieter Abbeel. PixelSNAIL: An improved autoregressive generative model. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 864–872. PMLR, 10–15 Jul 2018.
- 547 Minsik Cho, Keivan Alizadeh-Vahid, Saurabh Adya, and Mohammad Rastegari. Dkm: Differen 548 tiable k-means clustering layer for neural network compression. In *International Conference on Learning Representations*, 2021.
- Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Variable rate deep image compression with a conditional autoencoder. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3146–3154, 2019.
- Max Cohen, Guillaume Quispe, Sylvain Le Corff, Charles Ollion, and Eric Moulines. Diffusion bridges vector quantized variational autoencoders. In *International Conference on Machine Learning*, pp. 4141–4156. PMLR, 2022.
- <sup>557</sup> Ze Cui, Jing Wang, Bo Bai, Tiansheng Guo, and Yihui Feng. G-vae: A continuously variable rate
   <sup>558</sup> deep image compression framework. *arXiv preprint arXiv:2003.02012*, 2(3), 2020.
- Andrew M Dai and Quoc V Le. Semi-supervised sequence learning. Advances in neural information processing systems, 28, 2015.
- Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever.
   Jukebox: A generative model for music. *arXiv preprint arXiv:2005.00341*, 2020.
- Lyndon R Duong, Bohan Li, Cheng Chen, and Jingning Han. Multi-rate adaptive transform coding for video compression. In *ICASSP 2023-2023 IEEE International Conference on Acoustics*, *Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image
   synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recogni- tion*, pp. 12873–12883, 2021.
- William A Falcon. Pytorch lightning. *GitHub*, 3, 2019.

582

583

584

- Yunchao Gong, Liu Liu, Ming Yang, and Lubomir Bourdev. Compressing deep convolutional net works using vector quantization. *arXiv preprint arXiv:1412.6115*, 2014.
- 575 Robert Gray. Vector quantization. *IEEE Assp Magazine*, 1(2):4–29, 1984.
  576
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola.
  A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.
- Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and
   Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10696–10706, 2022.
  - Haohan Guo, Fenglong Xie, Frank K Soong, Xixin Wu, and Helen Meng. A multi-stage multicodebook vq-vae approach to high-performance neural tts. In *Proc. INTERSPEECH*, 2022.
- Liyong Guo, Xiaoyu Yang, Quandong Wang, Yuxiang Kong, Zengwei Yao, Fan Cui, Fangjun Kuang, Wei Kang, Long Lin, Mingshuang Luo, et al. Predicting multi-codebook vector quanti-zation indexes for knowledge distillation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
- 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. Advances in neural information processing systems, 30, 2017.
- <sup>593</sup> Iris Huijben, Matthijs Douze, Matthew Muckley, Ruud van Sloun, and Jakob Verbeek. Residual quantization with implicit neural codebooks. *arXiv preprint arXiv:2401.14732*, 2024.

| 594<br>595<br>596               | Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In <i>Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14</i> , pp. 694–711. Springer, 2016.   |
|---------------------------------|--|
| 597<br>598<br>599<br>600<br>601 | Nick Johnston, Damien Vincent, David Minnen, Michele Covell, Saurabh Singh, Troy Chinen, Sung Jin Hwang, Joel Shor, and George Toderici. Improved lossy image compression with priming and spatially adaptive bit rates for recurrent networks. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 4385–4393, 2018. |
| 602<br>603                      | Diederik P Kingma and Max Welling. Auto-encoding variational bayes. <i>arXiv preprint arXiv:1312.6114</i> , 2013.  |
| 604<br>605<br>606<br>607        | Jari Korhonen and Junyong You. Peak signal-to-noise ratio revisited: Is simple beautiful? In 2012 Fourth International Workshop on Quality of Multimedia Experience, pp. 37–38, 2012. doi: 10.1109/QoMEX.2012.6263880.   |
| 608<br>609<br>610               | Ashok K. Krishnamurthy, Stanley C. Ahalt, Douglas E. Melton, and Prakoon Chen. Neural networks for vector quantization of speech and images. <i>IEEE journal on selected areas in Communications</i> , 8(8):1449–1457, 1990.   |
| 611<br>612<br>613               | Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced re-<br>search). URL http://www.cs.toronto.edu/~kriz/cifar.html.  |
| 614<br>615<br>616<br>617        | Kundan Kumar, Rithesh Kumar, Thibault De Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre De Brebisson, Yoshua Bengio, and Aaron C Courville. Melgan: Generative adversarial networks for conditional waveform synthesis. <i>Advances in neural information processing systems</i> , 32, 2019.   |
| 618<br>619<br>620<br>621        | Alex M Lamb, Anirudh Goyal ALIAS PARTH GOYAL, Ying Zhang, Saizheng Zhang, Aaron C Courville, and Yoshua Bengio. Professor forcing: A new algorithm for training recurrent networks. <i>Advances in neural information processing systems</i> , 29, 2016.   |
| 622<br>623<br>624<br>625        | Adrian Łańcucki, Jan Chorowski, Guillaume Sanchez, Ricard Marxer, Nanxin Chen, Hans JGA Dolfing, Sameer Khurana, Tanel Alumäe, and Antoine Laurent. Robust training of vector quantized bottleneck models. In 2020 International Joint Conference on Neural Networks (IJCNN), pp. 1–7. IEEE, 2020.   |
| 626<br>627<br>628               | Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image generation using residual quantization. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 11523–11532, 2022.   |
| 629<br>630<br>631<br>632        | Lei Li, Tingting Liu, Chengyu Wang, Minghui Qiu, Cen Chen, Ming Gao, and Aoying Zhou. Resiz-<br>ing codebook of vector quantization without retraining. <i>Multimedia Systems</i> , 29(3):1499–1512, 2023.   |
| 633<br>634                      | Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.<br>In <i>Proceedings of International Conference on Computer Vision (ICCV)</i> , December 2015.  |
| 635<br>636<br>637<br>638        | Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id=Bkg6RiCqY7.  |
| 639<br>640<br>641<br>642        | Minh-Thang Luong, Quoc V. Le, Ilya Sutskever, Oriol Vinyals, and Lukasz Kaiser. Multi-task sequence to sequence learning. In Yoshua Bengio and Yann LeCun (eds.), 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings, 2016. URL http://arxiv.org/abs/1511.06114.          |
| 643<br>644<br>645               | May Malka, Shai Ginzach, and Nir Shlezinger. Learning multi-rate vector quantization for remote deep inference. In 2023 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW), pp. 1–5. IEEE, 2023.  |
| 647                             | Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. Finite scalar quantiza-<br>tion: Vq-vae made simple. <i>arXiv preprint arXiv:2309.15505</i> , 2023.  |

680

685

686

687

688

689

690

691 692

693

694

- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- Sangwoo Park, Osvaldo Simeone, and Joonhyuk Kang. Meta-learning to communicate: Fast end-to-end training for fading channels. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5075–5079. IEEE, 2020.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
   Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- 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*, pp. 8821–8831. PMLR, 2021.
- Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence level training with recurrent neural networks. In Yoshua Bengio and Yann LeCun (eds.), 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings, 2016. URL http://arxiv.org/abs/1511.06732.
- 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.
- 669 Danilo Jimenez Rezende and Fabio Viola. Taming vaes. arXiv preprint arXiv:1810.00597, 2018.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- Myungseo Song, Jinyoung Choi, and Bohyung Han. Variable-rate deep image compression through
   spatially-adaptive feature transform. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2380–2389, 2021.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks.
   Advances in neural information processing systems, 27, 2014.
- Yuhta Takida, Takashi Shibuya, Weihsiang Liao, Chieh-Hsin Lai, Junki Ohmura, Toshimitsu Uesaka, Naoki Murata, Shusuke Takahashi, Toshiyuki Kumakura, and Yuki Mitsufuji. Sq-vae:
  Variational bayes on discrete representation with self-annealed stochastic quantization. In *International Conference on Machine Learning*, pp. 20987–21012. PMLR, 2022.
  - Andros Tjandra, Berrak Sisman, Mingyang Zhang, Sakriani Sakti, Haizhou Li, and Satoshi Nakamura. Vqvae unsupervised unit discovery and multi-scale code2spec inverter for zerospeech challenge 2019. arXiv preprint arXiv:1905.11449, 2019.
  - Jonathan Tseng, Rodrigo Castellon, and Karen Liu. Edge: Editable dance generation from music. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 448–458, 2023.
  - Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. Advances in neural information processing systems, 30, 2017.
- Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), Proceedings of The 33rd International Conference on Machine Learning, volume 48 of Proceedings of Machine Learning Research, pp. 1747–1756, New York, New York, USA, 20–22 Jun 2016. PMLR. URL https://proceedings.mlr.press/v48/oord16.html.
- Benjamin Van Niekerk, Leanne Nortje, and Herman Kamper. Vector-quantized neural networks for acoustic unit discovery in the zerospeech 2020 challenge. *arXiv preprint arXiv:2005.09409*, 2020.

| 702<br>703<br>704        | Tung-Long Vuong, Trung Le, He Zhao, Chuanxia Zheng, Mehrtash Harandi, Jianfei Cai, and Dinh Phung. Vector quantized wasserstein auto-encoder. <i>arXiv preprint arXiv:2302.05917</i> , 2023.  |
|--------------------------|---|
| 705<br>706               | Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:<br>from error visibility to structural similarity. <i>IEEE transactions on image processing</i> , 13(4):600–<br>612, 2004   |
| 707<br>708               | Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent  |
| 709<br>710               | neural networks. Neural computation, 1(2):270–280, 1989.  |
| 711<br>712               | Will Williams, Sam Ringer, Tom Ash, David MacLeod, Jamie Dougherty, and John Hughes. Hier-<br>archical quantized autoencoders. <i>Advances in Neural Information Processing Systems</i> , 33:4524–  |
| 713<br>714               | 4555, 2020.   |
| 715<br>716               | <i>ceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pp. 6380–6387, 2020.   |
| 717<br>718<br>719        | Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong.<br>Codetalker: Speech-driven 3d facial animation with discrete motion prior. In <i>Proceedings of the</i><br><i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 12780–12790, 2023. |
| 720<br>721<br>722        | Jialong Xu, Tze-Yang Tung, Bo Ai, Wei Chen, Yuxuan Sun, and Deniz Deniz Gündüz. Deep joint source-channel coding for semantic communications. <i>IEEE Communications Magazine</i> , 61(11): 42–48, 2023.  |
| 723<br>724<br>725<br>726 | Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete diffusion model for text-to-sound generation. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 2023.  |
| 727<br>728<br>729        | Fei Yang, Luis Herranz, Joost Van De Weijer, José A Iglesias Guitián, Antonio M López, and Mikhail G Mozerov. Variable rate deep image compression with modulated autoencoder. <i>IEEE Signal Processing Letters</i> , 27:331–335, 2020.  |
| 730<br>731<br>732        | 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. In <i>International Conference on Learning Representations</i> , 2022.                                    |
| 733<br>734<br>735<br>736 | Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 586–595, 2018.                              |
| 737<br>738<br>739        | Chuanxia Zheng and Andrea Vedaldi. Online clustered codebook. In <i>Proceedings of the IEEE/CVF</i><br><i>International Conference on Computer Vision</i> , pp. 22798–22807, 2023.  |
| 740                      |   |
| 741                      |   |
| 742                      |   |
| 743                      |   |
| 744                      |   |
| 743                      |   |
| 740                      |   |
| 740                      |   |
| 740                      |   |
| 749                      |   |
| 751                      |   |
| 750                      |   |
| 752                      |   |
| 754                      |   |
| 755                      |   |
|                          |   |

#### A APPENDIX / SUPPLEMENTARY MATERIAL

#### A.1 VQ-VAE CODEBOOK UPDATES WITH EXPONENTIAL MOVING AVERAGES (EMA)

760 At training step t, the  $n_i$  encoder outputs  $\{f_{\phi}(x_1), f_{\phi}(x_2), ..., f_{\phi}(x_{n_i})\}$  from codebook  $e_i$  for the 761 mini-batch data  $\{x_1, x_2, ..., x_{n_i}\}$  are updated with count  $N_i^{(t)}$  and mean value  $m_i^{(t)}$  as follows:

$$N_{i}^{(t)} := \gamma \cdot N_{i}^{(t-1)} + (1-\gamma) \cdot n_{i}^{(t)}$$

$$m_{i}^{(t)} := \gamma \cdot m_{i}^{(t-1)} + (1-\gamma) \cdot \sum_{j}^{n_{i}^{(t)}} f_{\phi}(x_{j})^{(t)}$$

$$e_{i}^{(t)} := \frac{m_{i}^{(t)}}{N_{i}^{(t)}}$$
(5)

771

772

773 774

775

782

783 784 785

786 787

788 789

790

792

793

806

766 767

756

758

759

where a  $\gamma$  is a decay factor with a value between 0 and 1 (the default value  $\gamma = 0.99$  was used in all of our experiments). The count  $N_i^{(t)}$  represents the encoder hidden states that have  $e_i$  as it's nearest neighbor.  $N_i^{(0)}$  is initially set as zero.

#### A.2 CODEBOOK CLUSTERING OF MODEL-BASED RAQ-VAE

Given a set of the original codebook representations  $\mathbf{e} = \{e_i\}_{i=1}^{K}$ , we aim to partition the *K* code vectors into  $\widetilde{K}(\leq K)$  code vectors  $\widetilde{\mathbf{e}} = \{\widetilde{e}_i\}_{i=1}^{\widetilde{K}}$ . Each codebook vector resides in a *D*-dimensional Euclidean space. Using the codebook assignment function  $g(\cdot)$ , then  $g(e_i) = j$  means *i*-th given codebook assigned *j*-th clustered codebook. Our objective for codebook clustering is to minimize the discrepancy  $\mathcal{L}$  between the given codebook e and clustered codebook  $\widetilde{\mathbf{e}}$ :

$$\underset{\tilde{\mathbf{e}},g}{\arg\min} \mathcal{L}(\mathbf{e}; \tilde{\mathbf{e}}) = \underset{\tilde{\mathbf{e}},g}{\arg\min} \sum_{i=1}^{\tilde{K}} ||e_i - \tilde{e}_{g(e_i)}||$$
(6)

with necessary conditions

$$g(e_i) = \underset{j \in 1, 2, \dots, \tilde{K}}{\operatorname{arg\,min}} ||e_i - \tilde{e}_j||, \quad \tilde{e}_j = \frac{\sum_{i:g(e_i)=j} e_i}{N_j}$$
(7)

where  $N_i$  is the number of samples assigned to the codebook  $\tilde{e}_i$ .

#### 791 A.3 EXPERIMENT DETAILS

#### A.3.1 ARCHITECTURES AND HYPERPARAMETERS

794 The model architecture for this study is based on the conventional VQ-VAE framework outlined in the original VQ-VAE paper (Van Den Oord et al., 2017), and is implemented with reference to the 796 VQ-VAE-2 (Razavi et al., 2019) implementation repositories <sup>123</sup>. We are using the ConvResNets 797 from the repositories. These networks consist of convolutional layers, transpose convolutional layers 798 and ResBlocks. Experiments were conducted on two different computer setups: a server with 4 RTX 799 4090 GPUs and a machine with 2 RTX 3090 GPUs. PyTorch (Paszke et al., 2019), PyTorch Lightning (Falcon, 2019), and the AdamW (Loshchilov & Hutter, 2019) optimizer were used for model 800 implementation and training. Evaluation metrics such as the Structural Similarity Index (SSIM) and 801 the Frechet Inception Distance (rFID) were computed using implementations of pytorch-msssim<sup>4</sup> 802 and pytorch-fid<sup>5</sup>, respectively. The detailed model parameters are shown in Table 2. RAQ-VAEs 803 are constructed based on the described VQ-VAE parameters with additional consideration of each 804 parameter. 805

<sup>&</sup>lt;sup>1</sup>https://github.com/mattiasxu/VQVAE-2

<sup>&</sup>lt;sup>2</sup>https://github.com/rosinality/vq-vae-2-pytorch

<sup>&</sup>lt;sup>3</sup>https://github.com/EugenHotaj/pytorch-generative

<sup>&</sup>lt;sup>4</sup>https://github.com/VainF/pytorch-msssim

<sup>&</sup>lt;sup>5</sup>https://github.com/mseitzer/pytorch-fid

| 811 | Method                             | Parameter                               | CIFAR10           | CelebA             | ImageNet                  |
|-----|------------------------------------|---|-------------------|--------------------|---------------------------|
| 812 |                                    | Input size                              | 32×32×3           | 64×64×3            | 224×224×3                 |
| 813 |                                    | Latent layers                           | 8×8               | 16×16              | 56×56                     |
| 814 |                                    | Hidden units                            | 128               | 128                | 256                       |
| 915 |                                    | Residual units                          | 64                | 64                 | 128                       |
| 015 |                                    | # of ResBlock                           | $2 \\ 24 \\ 210$  | $2_{05}$ $0^{11}$  | $2 \\ 07 \\ 012$          |
| 816 | VQ-VAE (Van Den Oord et al., 2017) | Codebook dimension (D)                  | $2 \sim 2^{-1}$   | $2^{\circ} \sim 2$ | $2^{\circ} \sim 2$<br>128 |
| 817 |                                    | $\beta$ (Commit loss weight)            | 0.25              | 0.25               | 0.25                      |
| 818 |                                    | Weight decay in EMA ( $\gamma$ )        | 0.99              | 0.99               | 0.99                      |
| 819 |                                    | Batch size                              | 128               | 128                | 32                        |
| 000 |                                    | Optimizer                               | AdamW             | AdamW              | AdamW                     |
| 820 |                                    | Learning rate                           | 0.0005            | 0.0005             | 0.0005                    |
| 821 |                                    | Max. training steps                     | 195K              | 635.5K             | 2500K                     |
| 822 |                                    | Original codebook size $(K)$            | 64, 128           | 128, 256           | 512                       |
| 823 |                                    | Adapted codebook size $(\widetilde{K})$ | $2^4 \sim 2^{10}$ | $2^5 \sim 2^{11}$  | $2^6 \sim 2^{12}$         |
| 00/ | Model-based RAQ-VAE                | Max. DKM iteration                      | 200               | 200                | 200                       |
| 024 |                                    | Max. IKM iteration                      | 5000              | 5000               | 5000                      |
| 825 |                                    | au of softmax                           | 0.01              | 0.01               | 0.01                      |
| 826 |                                    | Original codebook size $(K)$            | 64, 128           | 128, 256           | 512                       |
| 827 |                                    | Adapted codebook size $(K)$             | $2^4 \sim 2^{10}$ | $2^5 \sim 2^{11}$  | $2^6 \sim 2^{12}$         |
| 828 | Data-driven RAO-VAF                | Max. Codebook size                      | 1024              | 2048               | 4096                      |
| 020 | Data-univen KAQ- VAE               | Min. Codebook size                      | 8<br>64           | 10<br>64           | 04<br>128                 |
| 829 |                                    | Hidden size (Seq2Seq)                   | 64                | 64                 | 128                       |
| 830 |                                    | # of recurrent layers (Seq2Seq)         | 2                 | 2                  | 2                         |
| 831 |                                    |   |                   |                    |                           |

Table 2: Architecture and hyperparameters of baseline VQ-VAE model and its RAQ-VAE model (Model-based RAQ-VAE and Data-driven RAQ-VAE)

### A.3.2 DATASETS AND PREPROCESSING

838 For the **CIFAR10** dataset, the training set is preprocessed using a combination of random cropping 839 and random horizontal flipping. Specifically, a random crop of size  $32 \times 32$  with padding of 4 using the 'reflect' padding mode is applied, followed by a random horizontal flip. The validation and 840 test sets are processed by converting the images to tensors without further augmentation. For the 841 **CelebA** dataset, the training set is preprocessed with a series of transformations. The images are 842 resized and center cropped to  $64 \times 64$ , normalized, and subjected to random horizontal flipping. A 843 similar preprocessing is applied to the validation set, while the test set is processed without augmen-844 tation. For the **ImageNet** dataset, the training set is preprocessed with a series of transformations. 845 The images are resized  $256 \times 256$  and center cropped to  $224 \times 224$ , normalized, and subjected to 846 random horizontal flipping. A similar preprocessing is applied to the validation set, while the test 847 set is processed without augmentation. These datasets are loaded into PyTorch using the provided 848 data modules, and the corresponding data loaders are configured with the specified batch sizes and 849 learning rate for efficient training (described in Table 2. The datasets are used as input for training, 850 validation, and testing of the VQ-VAE model.

851 852

810 81

832

833

834 835 836

837

#### A.3.3 MODEL COMPLEXITY

853 854

855 To provide a comprehensive understanding of the model complexity for the different datasets used in our experiments, we detail the number of parameters in the Encoder, Decoder, Quantizer, and 856 Seq2Seq components of the trained models in Table 3 and 4. The table summarizes the number of 857 model parameters counts for the CIFAR10 and CelebA datasets. 858

859 Moreover, we show the training/inference time in Table 5 and 6. The training and inference times 860 were measured for both model-based and data-driven RAQ-VAE methods. These results highlight the trade-offs between our two methods and their potential applications depending on resource avail-861 ability and performance requirements. Although the results show that the data-driven method has 862 a higher computational cost, we expect that the benefit of achieving the rate will provide better 863 flexibility and performance in different scenarios.

| Method  |         |         | # params  |         |       |
|---|---------|---------|-----------|---------|-------|
|   | Encoder | Decoder | Quantizer | Seq2Seq | Total |
| <b>VQ-VAE</b> ( $K = 1024$ )                              | 196.3K  | 262K    | 65.5 K    | -       | 525K  |
| <b>VQ-VAE</b> ( $K = 512$ )                               | 196.3K  | 262K    | 32.8K     | -       | 492K  |
| <b>VQ-VAE</b> ( $K = 256$ )                               | 196.3K  | 262K    | 16.4K     | -       | 476K  |
| <b>VQ-VAE</b> ( $K = 128$ )                               | 196.3K  | 262K    | 8.2K      | -       | 468K  |
| $\mathbf{VQ}\text{-}\mathbf{VAE}\ (K=64)$                 | 196.3K  | 262K    | 4.1K      | -       | 463K  |
| <b>VQ-VAE</b> ( $K = 32$ )                                | 196.3K  | 262K    | 2.0K      | -       | 461K  |
| VQ-VAE (K = 16)   | 196.3K  | 262K    | 1.0K      | -       | 460K  |
| <b>VQ-VAE</b> ( $K = 1024$ ) (randomly selected codebook) | 196.3K  | 262K    | 65.5 K    | -       | 525K  |
| <b>Data-driven RAQ-VAE</b> ( $K = 128$ )                  | 196.3K  | 262K    | 8.2K      | 263.7K  | 732K  |
| <b>Data-driven RAQ-VAE</b> $(K = 64)$                     | 196.3K  | 262K    | 4.1K      | 263.7K  | 728K  |
| <b>Model-based RAQ-VAE</b> ( $K = 128$ )                  | 196.3K  | 262K    | 8.2K      | -       | 468K  |
| <b>Model-based RAQ-VAE</b> $(K = 64)$                     | 196.3K  | 262K    | 4.1K      | -       | 463K  |

Table 3: Number of parameters for training our models on CIFAR10 dataset.

| Method  |         |         | # params  |         |       |
|---|---------|---------|-----------|---------|-------|
|   | Encoder | Decoder | Quantizer | Seq2Seq | Total |
| <b>VQ-VAE</b> ( $K = 2048$ )                              | 196.3K  | 262K    | 131K      | -       | 590K  |
| <b>VQ-VAE</b> ( $K = 1024$ )                              | 196.3K  | 262K    | 65.5 K    | -       | 525K  |
| <b>VQ-VAE</b> ( $K = 512$ )                               | 196.3K  | 262K    | 32.8K     | -       | 492K  |
| <b>VQ-VAE</b> ( $K = 256$ )                               | 196.3K  | 262K    | 16.4K     | -       | 476K  |
| <b>VQ-VAE</b> ( $K = 128$ )                               | 196.3K  | 262K    | 8.2K      | -       | 468K  |
| $\mathbf{VQ-VAE}(K=64)$                                   | 196.3K  | 262K    | 4.1K      | -       | 463K  |
| <b>VQ-VAE</b> ( $K = 32$ )                                | 196.3K  | 262K    | 2.0K      | -       | 461K  |
| <b>VQ-VAE</b> ( $K = 2048$ ) (randomly selected codebook) | 196.3K  | 262K    | 131K      | -       | 590K  |
| <b>Data-driven RAQ-VAE</b> $(K = 256))$                   | 196.3K  | 262K    | 16.4K     | 263.7K  | 740K  |
| <b>Data-driven RAQ-VAE</b> ( $K = 128$ )                  | 196.3K  | 262K    | 8.2K      | 263.7K  | 732K  |
| Model-based RAQ-VAE ( $K = 256$ )                         | 196.3K  | 262K    | 16.4K     | -       | 476K  |
| Model-based RAQ-VAE ( $K = 128$ )                         | 196.3K  | 262K    | 8.2K      | -       | 468K  |

Table 4: Number of parameters for training our models on CelebA dataset.

| Method                       | K    | Training time per epoch (s) | # params |
|------------------------------|------|-----------------------------|----------|
| VQ-VAE / Model-based RAQ-VAE | 64   | $18.09\pm0.256$             | 463K     |
| VQ-VAE / Model-based RAQ-VAE | 256  | $18.43\pm0.1$               | 476K     |
| VQ-VAE / Model-based RAQ-VAE | 1024 | $21.64\pm0.11$              | 525K     |
| Data-driven RAQ-VAE          | 256  | $514.97 \pm 08.17$          | 740K     |

Table 5: Training time per epoch on on CelebA train set using a Nvidia RTX 3090 GPU.

## A.4 ADDITIONAL EXPERIMENTS

913 A.4.1 REDUCING THE RATE 

As analyzed in Section 5.1, data-driven RAQ-VAE generally outperforms model-based RAQ-VAE,
 but some rate-reduction results on CIFAR10 show that model-based RAQ-VAE performs much more
 stably than in the codebook increasing task. This indicates that simply clustering codebook vectors, without additional neural models like Seq2Seq, can achieve remarkable performance.

|       | Method   |   | $\widetilde{K}$ I   | nference ti  | me ner enoch (s)  |
|-------|--|---|---|--|---|
|       | memou  |   | 11 I.   | inci che u   | me per epoen (s)  |
|       | <b>VQ-VAE</b> $(K = 64)$   |   | -   | 1.8  | $6 \pm 0.10$  |
|       | <b>VQ-VAE</b> ( $K = 256$ )  |   | -   | 1.9  | $1 \pm 0.12$  |
|       | <b>VQ-VAE</b> ( $K = 1024$ )   |   | -   | 1.8  | $6 \pm 0.09$  |
|       | Model-based RAQ-VAE ( $K =$  | = 256)  | 64  | 1.9  | $8 \pm 0.09$  |
|       | <b>Data-driven RAQ-VAE</b> ( $K =$   | 256)  | 64  | 3.0  | $5 \pm 0.11$  |
|       | Model-based RAQ-VAE ( $K =$  | = 256)  | 1024  | 70.9   | $1 \pm 11.82$   |
|       | <b>Data-driven RAQ-VAE</b> ( $K =$   | 256)  | 1024  | 33.2   | $21 \pm 0.27$   |
| Table | e 6: Inference time per epoch o  | n Celel   | bA test set   | using a N  | vidia RTX 309   |
| Table | e 6: Inference time per epoch o  | $\widetilde{K}$   | bA test set   | Tusing a N   | = 1024)   |
| Table | e 6: Inference time per epoch o<br>Method  | on Celel $\widetilde{K}$  | bA test set<br>CII<br>PSNR↑                                       | using a N<br>FAR10 ( <i>K</i><br>rFID↓   | lvidia RTX 309<br>= 1024)<br>Perplexity ↑                             |
| Table | e 6: Inference time per epoch o Method VQ-VAE (baseline model)                                 | on Celel $\widetilde{K}$  | DA test set<br>CIII<br>PSNR↑<br>  25.48                           | x using a N<br>FAR10 ( <i>K</i><br>rFID ↓<br>51.90   | lvidia RTX 309<br>= $1024$ )<br>Perplexity $\uparrow$<br>708.60       |
| Table | e 6: Inference time per epoch o Method VQ-VAE (baseline model)                                 | $\widetilde{K}$ $\overline{\tilde{K}}$ $\overline{512}$                                   | bA test set<br>CIII<br>PSNR↑<br>25.48<br>24.35                    | For the second | lvidia RTX 309<br>$= 1024)$ Perplexity ↑ $\overline{708.60}$ $289.29$ |
| Table | e 6: Inference time per epoch o Method VQ-VAE (baseline model) VQ-VAE (random select)          | $ \begin{array}{c c} \widetilde{K} \\ \widetilde{K} \\ \hline \\ 512 \\ 256 \end{array} $ | bA test set<br>PSNR↑<br>25.48<br>24.35<br>22.81                   | FAR10 ( $K$<br>rFID $\downarrow$<br>51.90<br>63.67<br>78.00  |   |
| Table | e 6: Inference time per epoch o<br>Method<br>VQ-VAE (baseline model)<br>VQ-VAE (random select) | n Celel<br><i>K</i> - 512 256 128   | bA test set<br>PSNR↑<br>25.48<br>24.35<br>22.81<br>20.87          | Example a second strain from the second strain |   |
| Table | e 6: Inference time per epoch o Method VQ-VAE (baseline model) VQ-VAE (random select)          | K           -           512           256           128           512                     | bA test set<br>PSNR↑<br>25.48<br>24.35<br>22.81<br>20.87<br>24.62 | Fusing a N         FAR10 ( $K$ rFID $\downarrow$ 51.90         63.67         78.00         93.57 <b>55.78</b>  | $= 1024)$ Perplexity $\uparrow$ 708.60 289.29 111.77 48.87 285.68     |

| vQ-vAE (baseline model)        | -                          | 25.48                                   | 51.90                                    | /08.60                                    |  |
|--------------------------------|----------------------------|---|--|---|--|
| VQ-VAE (random select)         | 512<br>256                 | 24.35<br>22.81                          | 63.67<br>78.00                           | <b>289.29</b><br>111.77                   |  |
|                                | 128                        | 20.87                                   | 93.57                                    | 48.87                                     |  |
| Model-based RAQ-VAE            | 512<br>256<br>128          | 24.62<br>23.81<br>23.07                 | 55.78<br>62.53<br>69.45                  | 285.68<br>134.54<br>73.17                 |  |
| Method                         | $\widetilde{K}$            | <b>CelebA</b> ( $K = 2048$ )            |  |   |  |
|                                | 11                         | PSNR ↑                                  | rFID↓                                    | Perplexity ↑                              |  |
| <b>VQ-VAE</b> (baseline model) | -                          | 28.26                                   | 22.89                                    | 273.47                                    |  |
|                                |                            |   |  |   |  |
|                                | 1024                       | 24.02                                   | 38.92                                    | 103.50                                    |  |
| <b>VQ-VAE</b> (random select)  | 1024<br>512                | 24.02<br>18.99                          | 38.92<br>71.64                           | <b>103.50</b><br>49.59                    |  |
| VQ-VAE (random select)         | 1024<br>512<br>256         | 24.02<br>18.99<br>23.54                 | 38.92<br>71.64<br>115.12                 | <b>103.50</b><br>49.59<br>27.86           |  |
| VQ-VAE (random select)         | 1024<br>512<br>256<br>1024 | 24.02<br>18.99<br>23.54<br><b>26.40</b> | 38.92<br>71.64<br>115.12<br><b>31.37</b> | <b>103.50</b><br>49.59<br>27.86<br>102.36 |  |

Table 7: Reconstuction performances for **rate-reduction task** according to adapted codebook size  $\widetilde{K}$ . The distortion (PSNR), perceptual similarity (rFID), and codebook usability (perplexity) are evaluated using the test set on CIFAR-10 an CelebA. Higher values are better for PSNR, and perplexity, while lower values are better for rFID.

In Table 7, the performance via codebook clustering was evaluated with different original/adapted codebook sizes K: 1024 / K: 512, 256, 128 on CIFAR10 and K: 2048 / K: 1024, 512, 256, 128 on CelebA. The conventional VQ-VAE preserved as many codebooks in the original codebook as in the adapted codebook, while randomly codebook-selected VQ-VAE results remained meaningless. Model-based RAQ-VAE adopted this baseline VQ-VAE model and performed clustering on the adapted codebook. Model-based RAQ-VAE shows a substantial performance difference in terms of reconstructed image distortion and codebook usage compared to randomly codebook-selected VQ-VAE. Even when evaluating absolute performance, it is intuitive that online codebook representation via model-based RAQ-VAE provides some performance guarantees. 

963 A.4.2 INCREASING THE RATE

In our proposed RAQ-VAE scenario, increasing the codebook size beyond the base size is a more demanding and crucial task than reducing it. The crucial step in building data-driven RAQ-VAE is to achieve higher rates from a fixed model architecture and compression rate, ensuring usability. Therefore, the codebook increasing task was the main challenge. The Seq2Seq decoding algorithm based on cross-forcing is designed with this intention.

In Figure 2, the codebook generation performance was evaluated with different original/adapted codebook sizes K: 64, 128 /  $\tilde{K}$ : 64, 128, 256, 512, 1024 on CIFAR10 and K: 128, 256 /  $\tilde{K}$ : 128, 256, 512, 1024, 2048 on CelebA datasets. As discussed in Section 5.1, data-driven RAQ-VAE out-



Figure 4: **Reconstruction performance** at different rates (adapted codebook sizes) evaluated on CelebA ( $64 \times 64$ ) test set. In the graph, the black VQ-VAE-2s (Razavi et al., 2019) are separate models trained on each codebook size, while the RAQ-VAEs are one model per line.

performs model-based RAQ-VAE in the rate-increasing task and partially outperforms conventional VQ-VAE trained on the same codebook size  $(K = \tilde{K})$ . This effect is particularly pronounced on CelebA.

However, increasing the difference between the original and adapted codebook sizes leads to a degradation of RAQ-VAE performance. This effect is more dramatic for model-based RAQ-VAE due to
its algorithmic limitations, making its performance less stable at high rates. Improving the performance of model-based RAQ-VAE, such as modifying the initialization of the codebook vector,
remains a limitation.

#### 996 A.4.3 Additional Quantitative Results

<sup>997</sup> <sup>998</sup> In Table 8 and 9, we present additional quantitative results for the reconstruction on CIFAR10 and <sup>999</sup> CelebA datasets. The error indicates a 95.45% confidence interval based on 4 runs with different <sup>1000</sup> training seeds. Figure 4 shows the reconstruction performance using VQ-VAE-2 as the baseline <sup>1001</sup> model. The results demonstrate that the data-driven RAQ-VAE model significantly outperforms the <sup>1001</sup> original VQ-VAE-2 across multiple rates on the CelebA ( $64 \times 64$ )dataset.

1002 1003

1006

983

984

985 986

A.4.4 ADDITIONAL QUALITATIVE RESULTS

1005 In Figure 5, we present additional qualitative results for the reconstruction on ImageNet dataset.

1007 A.4.5 EFFECTIVENESS OF Cross-forcing

We conducted additional experiments to demonstrate the effectiveness of the cross-forcing strategy, following a concern about the learning stability of this approach. In Table 10, the results compare the reconstruction performance of the data-driven RAQ-VAE (K = 128) with and without crossforcing on the CelebA test dataset.

Our experiments demonstrate that the cross-forcing strategy is no less stable than the data-driven approach without cross-forcing when trained with the same four seeds. Furthermore, the performance improvement from cross-forcing becomes significant, particularly when operating with a codebook size equal to or larger than twice the original codebook size, as intended. This is in line with the general goal in machine learning of training models with a smaller original codebook size while still achieving better reconstruction performance at higher bitrates. It consistently improves performance metrics such as MSE, PSNR and rFID as the codebook size increases, making it an effective approach for tasks that demand higher bitrates without compromising model efficiency.

- 1020
- 1021
- 1023
- 1024
- 1025

Bit Rate

| 4 | U | J | 0 |  |
|---|---|---|---|--|
| 1 | 0 | 3 | 7 |  |

| Method  | Bit Rate Codebook Usability   |  | Distortion  | Perceptual Similarity   |  |  |
|---|---|--|---|---|--|--|
|   | $\widetilde{K}(\mathrm{bpp})$   | Usage  | Perplexity  | PSNR  | rFID   | SSIM   |
| $\mathbf{VQ-VAE} \ (K = \widetilde{K})$                       | 1024 (0.625)  | 972.66±2.97  | 708.60±7.04   | 25.48±0.02  | 51.90±0.51   | 0.8648±0.0005  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$                           | 512 (0.5625)  | $507.52 {\pm} 0.51$  | 377.08±5.92   | $24.94{\pm}0.01$  | 56.65±0.91   | $0.8490 {\pm} 0.0003$  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$                           | 256 (0.5)   | 256±0  | 204.43±4.36   | 24.43±0.02  | 61.40±0.78   | 0.8310±0.0006  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$                           | 128 (0.4375)  | $128\pm0$  | $106.44{\pm}1.54$   | $23.85{\pm}0.01$  | 66.70±1.12   | 0.8096±0.0009  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$                           | 64 (0.375)  | 64±0   | 55.64±0.27  | 23.24±0.01  | 74.00±1.64   | 0.7849±0.0009  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$                           | 32 (0.3125)   | 32±0   | 29.25±0.13  | $22.53{\pm}0.02$  | 81.68±1.01   | 0.7545±0.0009  |
| $\mathbf{VQ}\text{-}\mathbf{VAE}\left(K=\widetilde{K}\right)$ | 16 (0.25)   | 16±0   | 15.01±0.21  | 21.76±0.01  | 89.75±0.83   | 0.7156±0.0024  |
| <b>VQ-VAE</b><br>(K = 1024)<br>(random select)                | 1024 (0.625)<br>512 (0.5625)<br>256 (0.5)<br>128 (0.4375)<br>64 (0.375)<br>32 (0.3125)  | $\begin{array}{c} 972.66{\pm}2.97\\ 498.38{\pm}1.85\\ 253.01{\pm}0.66\\ 127.34{\pm}0.33\\ 64{\pm}0\\ 32{\pm}0 \end{array}$                   | $708.60\pm7.04\\289.29\pm16.67\\111.77\pm21.53\\48.87\pm11.31\\24.31\pm5.26\\13.50\pm1.45$  | $\begin{array}{c} 25.48 {\pm} 0.02 \\ 24.35 {\pm} 0.11 \\ 22.81 {\pm} 0.38 \\ 20.87 {\pm} 0.73 \\ 19.46 {\pm} 0.98 \\ 17.76 {\pm} 1.12 \end{array}$                     | $51.90\pm0.51\\63.67\pm2.49\\78.00\pm5.07\\93.57\pm9.87\\109.90\pm14.20\\126.57\pm15.89$   | $\begin{array}{c} 0.8648 {\pm} 0.0005 \\ 0.8305 {\pm} 0.0056 \\ 0.7822 {\pm} 0.0100 \\ 0.7254 {\pm} 0.0235 \\ 0.6720 {\pm} 0.0309 \\ 0.6102 {\pm} 0.0350 \end{array}$                        |
| <b>Data-driven</b><br><b>RAQ-VAE</b><br>(K = 128)             | $\begin{array}{c} 1024(0.625)\\ 512\ (0.5625)\\ 256\ (0.5)\\ 128\ (0.4375)\\ 64\ (0.375)\\ 32\ (0.3125)\\ 16\ (0.25) \end{array}$ | $\begin{array}{r} 971.21 {\pm}4.14 \\ 503.48 {\pm}0.75 \\ 253.45 {\pm}0.50 \\ 128 {\pm}0 \\ 64 {\pm}0 \\ 32 {\pm}0 \\ 16 {\pm}0 \end{array}$ | $724.91{\pm}15.34\\380.02{\pm}6.82\\194.27{\pm}2.38\\109.65{\pm}3.50\\55.64{\pm}0.27\\29.50{\pm}0.21\\15.11{\pm}0.67$                                       | $\begin{array}{c} 24.85{\pm}0.02\\ 24.57{\pm}0.02\\ 24.12{\pm}0.02\\ 23.71{\pm}0.01\\ 23.08{\pm}0.02\\ 21.76{\pm}0.06\\ 20.79{\pm}0.18 \end{array}$                     | $\begin{array}{c} 57.03 \pm 1.34 \\ 59.36 \pm 0.62 \\ 62.18 \pm 1.09 \\ 66.89 \pm 1.07 \\ 71.84 \pm 0.31 \\ 82.85 \pm 0.87 \\ 104.86 \pm 5.91 \end{array}$             | $\begin{array}{c} 0.8420 {\pm} 0.0008 \\ 0.8326 {\pm} 0.0008 \\ 0.8193 {\pm} 0.0009 \\ 0.8071 {\pm} 0.0014 \\ 0.7855 {\pm} 0.0005 \\ 0.7384 {\pm} 0.0007 \\ 0.6918 {\pm} 0.0084 \end{array}$ |
| Model-based<br>RAQ-VAE<br>(K = 128)                           | 1024 (0.625)<br>512 (0.5625)<br>256 (0.5)<br>128 (0.4375)<br>64 (0.375)<br>32 (0.3125)<br>16 (0.25)                               | $744.36{\pm}18.74\\430.06{\pm}11.58\\244.61{\pm}3.13\\128{\pm}0\\64{\pm}0\\32{\pm}0\\16{\pm}0$   | $395.23\pm2.77$<br>$256.23\pm7.50$<br>$185.02\pm3.31$<br>$106.44\pm1.54$<br>$49.55\pm1.29$<br>$25.65\pm0.76$<br>$13.79\pm0.06$                              | $\begin{array}{c} 24.15{\pm}0.03\\ 24.04{\pm}0.03\\ 23.93{\pm}0.01\\ 23.85{\pm}0.01\\ 22.85{\pm}0.55\\ 21.88{\pm}0.75\\ 20.89{\pm}0.04 \end{array}$                     | $\begin{array}{c} 63.88 {\pm} 126 \\ 64.74 {\pm} 0.96 \\ 65.65 {\pm} 1.12 \\ 66.70 {\pm} 1.12 \\ 72.61 {\pm} 0.77 \\ 82.12 {\pm} 1.74 \\ 95.03 {\pm} 0.34 \end{array}$ | $\begin{array}{c} 0.8213 {\pm} 0.0014 \\ 0.8177 {\pm} 0.0012 \\ 0.8139 {\pm} 0.0010 \\ 0.8096 {\pm} 0.0009 \\ 0.7780 {\pm} 0.0013 \\ 0.7405 {\pm} 0.0046 \\ 0.6972 {\pm} 0.0010 \end{array}$ |
| <b>Data-driven</b><br><b>RAQ-VAE</b><br>(K = 64)              | 1024 (0.625)<br>512 (0.5625)<br>256 (0.5)<br>128 (0.4375)<br>64 (0.375)<br>32 (0.3125)<br>16 (0.25)                               | $\begin{array}{c} 972.14{\pm}6.49\\ 506.38{\pm}1.23\\ 255.52{\pm}0.48\\ 128{\pm}0\\ 64{\pm}0\\ 32{\pm}0\\ 16{\pm}0\\ \end{array}$            | $\begin{array}{c} 725.55{\pm}10.90\\ 382.43{\pm}10.58\\ 196.17{\pm}9.95\\ 109.65{\pm}3.50\\ 56.31{\pm}0.46\\ 29.62{\pm}0.66\\ 15.11{\pm}0.67 \end{array}$   | $\begin{array}{c} 25.04{\pm}0.01\\ 24.70{\pm}0.02\\ 24.25{\pm}0.02\\ 23.71{\pm}0.01\\ 23.23{\pm}0.01\\ 21.84{\pm}0.09\\ 20.79{\pm}0.18 \end{array}$                     | $55.34\pm1.48$<br>$57.91\pm1.42$<br>$61.96\pm1.00$<br>$66.89\pm1.07$<br>$71.17\pm1.17$<br>$90.04\pm1.44$<br>$104.86\pm5.91$  | $\begin{array}{c} 0.8487 {\pm} 0.0012 \\ 0.8387 {\pm} 0.0011 \\ 0.8245 {\pm} 0.0012 \\ 0.8071 {\pm} 0.0014 \\ 0.7897 {\pm} 0.0013 \\ 0.7350 {\pm} 0.0038 \\ 0.6918 {\pm} 0.0084 \end{array}$ |
| Model-based<br>RAQ-VAE<br>(K = 64)                            | 1024 (0.625)<br>512 (0.5625)<br>256 (0.5)<br>128 (0.4375)<br>64 (0.375)<br>32 (0.3125)<br>16 (0.25)                               | $706.20 \pm 115.18 \\ 428.39 \pm 12.29 \\ 233.75 \pm 4.63 \\ 125.07 \pm 1.58 \\ 64 \pm 0 \\ 32 \pm 0 \\ 16 \pm 0$                            | $\begin{array}{c} 345.50{\pm}107.06\\ 231.41{\pm}14.64\\ 140.19{\pm}2.82\\ 101.16{\pm}16.04\\ 55.64{\pm}0.27\\ 26.21{\pm}0.95\\ 13.59{\pm}0.85 \end{array}$ | $\begin{array}{c} 23.65 {\pm} 0.13 \\ 23.55 {\pm} 0.04 \\ 23.39 {\pm} 0.05 \\ 23.32 {\pm} 0.05 \\ 23.24 {\pm} 0.01 \\ 22.07 {\pm} 0.13 \\ 20.88 {\pm} 0.23 \end{array}$ | $\begin{array}{c} 70.30{\pm}2.02\\ 71.01{\pm}1.38\\ 71.72{\pm}1.43\\ 72.68{\pm}1.47\\ 74.00{\pm}1.64\\ 81.61{\pm}2.26\\ 92.84{\pm}3.30\\ \end{array}$                  | $\begin{array}{c} 0.8013 {\pm} 0.0051 \\ 0.7988 {\pm} 0.0005 \\ 0.7935 {\pm} 0.0012 \\ 0.7901 {\pm} 0.0008 \\ 0.7849 {\pm} 0.0009 \\ 0.7569 {\pm} 0.0014 \\ 0.7004 {\pm} 0.0063 \end{array}$ |

Codebook Usability

Distortion

Perceptual Similarity

Table 8: Reconstruction performance on CIFAR10 dataset. The 95.45% confidence interval is provided based on 4 runs with different training seeds.

1117

| Method  | Bit Rate  | Codebook Usability   |   | Distortion  | Perceptual Similarity  |  |
|---|---|--|---|---|--|--|
|   | $\widetilde{K}(\mathrm{bpp})$   | Usage  | Perplexity  | PSNR  | rFID   | SSIM   |
| <b>VQ-VAE</b> $(K = \widetilde{K})$               | 2048 (0.6875)   | 779.07±8.35  | 273.47±6.86   | 28.26±0.03  | 22.89±0.71   | 0.8890±0.0   |
| <b>VQ-VAE</b> $(K = \widetilde{K})$               | 1024 (0.625)  | 456.86±3.53  | 160.35±2.73   | 27.73±0.05  | 26.67±1.43   | 0.8763±0.0   |
| <b>VQ-VAE</b> $(K = \widetilde{K})$               | 512 (0.5625)  | 259.59±3.99  | 95.09±1.28  | 27.11±0.01  | 29.77±0.95   | 0.8636±0.0   |
| <b>VQ-VAE</b> $(K = \widetilde{K})$               | 256 (0.5)   | 144.44±2.49  | 57.86±0.91  | 26.46±0.03  | 31.53±1.01   | 0.8481±0.  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$               | 128 (0.4375)  | 80.26±0.99   | 34.98±0.39  | 25.72±0.04  | 36.25±0.98   | 0.8279±0.  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$               | 64 (0.375)  | 44.94±1.03   | 20.04±0.37  | 24.78±0.03  | 41.22±0.77   | 0.7986±0.  |
| <b>VQ-VAE</b> $(K = \widetilde{K})$               | 32 (0.3125)   | 25.48±0.69   | 12.69±0.31  | 23.76±0.06  | 46.56±1.97   | 0.7660±0.  |
| V <b>Q-VAE</b><br>K = 2048)<br>random select)     | 2048 (0.625)<br>1024 (0.5625)<br>512 (0.5)<br>256 (0.4375)  | $779.07 \pm 8.35$<br>$384.31 \pm 6.76$<br>$210.69 \pm 9.23$<br>$115.33 \pm 7.73$   | $\begin{array}{c} 273.47{\pm}6.86\\ 103.50{\pm}3.28\\ 49.59{\pm}4.54\\ 27.86{\pm}3.39 \end{array}$  | $\begin{array}{c} 28.26{\pm}0.03\\ 24.02{\pm}1.10\\ 18.99{\pm}1.40\\ 16.33{\pm}0.61\end{array}$   | $\begin{array}{c} 22.89{\pm}0.71\\ 38.92{\pm}3.27\\ 71.64{\pm}8.27\\ 115.12{\pm}11.93 \end{array}$   | 0.8890±0.<br>0.7963±0.<br>0.7037±0.<br>0.6353±0.   |
| <b>Data-driven</b><br><b>RAQ-VAE</b><br>(K = 256) | 2048 (0.625)<br>1024 (0.5625)<br>512 (0.5)<br>256 (0.4375)<br>128 (0.375)<br>64 (0.3125)<br>32 (0.25) | $\begin{array}{c} 885.53{\pm}6.76\\ 490.86{\pm}4.98\\ 275.84{\pm}1.72\\ 144.79{\pm}1.21\\ 80.21{\pm}4.27\\ 42.93{\pm}1.61\\ 22.76{\pm}1.57\\ \end{array}$    | $\begin{array}{r} 347.99 {\pm} 5.17 \\ 187.33 {\pm} 10.37 \\ 104.61 {\pm} 5.00 \\ 52.63 {\pm} 0.28 \\ 32.23 {\pm} 3.87 \\ 20.85 {\pm} 1.22 \\ 12.32 {\pm} 0.91 \end{array}$ | $\begin{array}{c} 27.96 {\pm} 0.14 \\ 27.51 {\pm} 0.13 \\ 26.95 {\pm} 0.086 \\ 26.29 {\pm} 0.054 \\ 25.13 {\pm} 0.26 \\ 24.09 {\pm} 0.21 \\ 22.62 {\pm} 0.27 \end{array}$ | $\begin{array}{c} 23.02{\pm}0.33\\ 25.08{\pm}0.23\\ 27.96{\pm}0.49\\ 32.34{\pm}0.86\\ 39.67{\pm}2.29\\ 51.57{\pm}6.66\\ 69.65{\pm}9.49\end{array}$   | 0.8858±0.<br>0.8758±0.<br>0.8637±0.<br>0.8463±0.<br>0.8162±0.<br>0.7912±0.<br>0.7479±0.  |
| Model-based<br>RAQ-VAE<br>(K = 256)               | 2048 (0.625)<br>1024 (0.5625)<br>512 (0.5)<br>256 (0.4375)<br>128 (0.375)<br>64 (0.3125)<br>32 (0.25) | $\begin{array}{c} 704.17{\pm}108.04\\ 460.77{\pm}26.98\\ 279.53{\pm}9.48\\ 144.44{\pm}2.49\\ 75.31{\pm}3.09\\ 41.66{\pm}1.22\\ 22.96{\pm}0.90\\ \end{array}$ | $\begin{array}{c} 117.53 \pm 33.57 \\ 134.48 \pm 11.26 \\ 100.64 \pm 08.94 \\ 57.86 \pm 0.91 \\ 25.05 \pm 1.95 \\ 14.73 \pm 0.56 \\ 3 \\ 10.16 \pm 0.95 \end{array}$        | $\begin{array}{c} 26.54{\pm}0.10\\ 26.59{\pm}0.06\\ 26.40{\pm}0.08\\ 26.46{\pm}0.03\\ 24.44{\pm}0.25\\ 22.85{\pm}0.36\\ 21.81{\pm}0.45\\ \end{array}$                     | $\begin{array}{c} 30.34{\pm}1.39\\ 30.49{\pm}1.10\\ 30.95{\pm}0.98\\ 31.53{\pm}1.01\\ 38.95{\pm}2.91\\ 48.96{\pm}1.13\\ 62.46{\pm}0.00 \end{array}$  | $\begin{array}{c} 0.8507 {\pm} 0.\\ 0.8509 {\pm} 0.\\ 0.8488 {\pm} 0.\\ 0.8481 {\pm} 0.\\ 0.7890 {\pm} 0.\\ 0.7391 {\pm} 0.\\ 0.7077 {\pm} 0. \end{array}$ |
| <b>Data-driven</b><br><b>RAQ-VAE</b><br>(K = 128) | 2048 (0.625)<br>1024 (0.5625)<br>512 (0.5)<br>256 (0.4375)<br>128 (0.375)<br>64 (0.3125)<br>32 (0.25) | $\begin{array}{c} 891.13{\pm}7.11\\ 490.15{\pm}14.39\\ 272.60{\pm}2.08\\ 152.65{\pm}2.45\\ 79.17{\pm}0.93\\ 42.71{\pm}1.66\\ 22.42{\pm}1.92 \end{array}$     | $345.25\pm5.15$<br>$176.71\pm6.19$<br>$96.87\pm2.68$<br>$60.90\pm2.18$<br>$31.36\pm0.77$<br>$19.78\pm2.31$<br>$11.43\pm2.14$  | $\begin{array}{c} 27.91 {\pm} 0.04 \\ 27.47 {\pm} 0.07 \\ 26.90 {\pm} 0.05 \\ 26.18 {\pm} 0.18 \\ 25.53 {\pm} 0.06 \\ 24.10 {\pm} 0.11 \\ 22.74 {\pm} 0.54 \end{array}$   | $\begin{array}{c} 22.64{\pm}0.76\\ 24.67{\pm}0.80\\ 26.90{\pm}0.04\\ 30.81{\pm}1.59\\ 36.30{\pm}1.12\\ 47.63{\pm}5.82\\ 62.39{\pm}3.76\end{array}$   | $\begin{array}{c} 0.8810\pm 0.\\ 0.8710\pm 0.\\ 0.8589\pm 0.\\ 0.8391\pm 0.\\ 0.8209\pm 0.\\ 0.7892\pm 0.\\ 0.7414\pm 0.\\ \end{array}$                    |
| Model-based<br>RAQ-VAE<br>(K = 128)               | 2048 (0.625)<br>1024 (0.5625)<br>512 (0.5)<br>256 (0.4375)<br>128 (0.375)<br>64 (0.3125)<br>32 (0.25) | $\begin{array}{r} 350.02{\pm}100.57\\ 432.15{\pm}45.80\\ 262.78{\pm}29.47\\ 153.16{\pm}5.46\\ 80.26{\pm}0.99\\ 41.88{\pm}0.72\\ 23.31{\pm}0.89 \end{array}$  | $\begin{array}{c} 64.87{\pm}21.22\\ 102.79{\pm}17.34\\ 75.63{\pm}12.04\\ 53.22{\pm}4.62\\ 34.98{\pm}0.39\\ 16.70{\pm}0.43\\ 9.56{\pm}0.77\\ \end{array}$                    | $\begin{array}{c} 22.77{\pm}0.78\\ 25.57{\pm}0.19\\ 25.50{\pm}0.29\\ 25.42{\pm}0.28\\ 25.72{\pm}0.04\\ 23.63{\pm}0.16\\ 21.64{\pm}0.13\\ \end{array}$                     | $\begin{array}{c} 52.37{\pm}10.94\\ 35.62{\pm}1.46\\ 36.82{\pm}0.73\\ 36.78{\pm}1.27\\ 36.25{\pm}0.98\\ 47.09{\pm}4.09\\ 64.85{\pm}6.92 \end{array}$ | 0.7463±0.<br>0.8296±0.<br>0.8265±0.<br>0.8285±0.<br>0.8279±0.<br>0.7736±0.<br>0.7037±0.  |

Table 9: Reconstruction performance on CelebA dataset. The 95.45% confidence interval is provided based on 4 runs with different training seeds.

| 18 |                   |                 |                     |                    |                    |                       |
|----|-------------------|-----------------|---------------------|--------------------|--------------------|-----------------------|
| 19 | Method            | $\widetilde{K}$ | $MSE\downarrow$     | PSNR $\uparrow$    | rFID $\downarrow$  | SSIM $\uparrow$       |
| 20 |                   | 2048 (†)        | 1.618±0.016         | 27.91±0.04         | 22.64±0.76         | 0.8810±0.0013         |
| 21 |                   | 1024 (†)        | $1.794{\pm}0.027$   | $27.47 {\pm} 0.07$ | $24.67 {\pm} 0.80$ | $0.8710{\pm}0.0016$   |
| 22 |                   | 512 (†)         | $2.042{\pm}0.021$   | $26.90{\pm}0.05$   | $26.90{\pm}0.04$   | $0.8589 {\pm} 0.0044$ |
| 3  | w/ cross-forcing  | 256 (†)         | $2.412{\pm}0.101$   | $26.18{\pm}0.18$   | 30.81±1.59         | $0.8391 {\pm} 0.0125$ |
| 4  |                   | 128 (-)         | $2.801 {\pm} 0.039$ | $25.53 {\pm} 0.06$ | $36.30{\pm}1.12$   | $0.8209 {\pm} 0.0072$ |
| -  |                   | 64 (↓)          | $3.895 {\pm} 0.095$ | $24.10 \pm 0.11$   | $47.63 {\pm} 5.82$ | $0.7892{\pm}0.0067$   |
|    |                   | 32 (↓)          | $5.357 {\pm} 0.630$ | $22.74{\pm}0.54$   | 62.39±3.76         | $0.7414{\pm}0.0304$   |
| ·  |                   | 2048 (†)        | $1.661 {\pm} 0.056$ | $27.80{\pm}0.14$   | 23.58±0.26         | $0.8789 {\pm} 0.0030$ |
|    |                   | 1024 (†)        | $1.815 {\pm} 0.050$ | $27.42 \pm 0.12$   | $25.46 {\pm} 0.26$ | $0.8705 {\pm} 0.0024$ |
|    |                   | 512 (†)         | $2.068 {\pm} 0.059$ | $26.85 {\pm} 0.12$ | $27.81 \pm 0.42$   | $0.8567 {\pm} 0.0046$ |
|    | w/o cross-forcing | 256 (†)         | $2.449 {\pm} 0.052$ | $26.12 \pm 0.09$   | $32.32{\pm}1.20$   | $0.8407{\pm}0.0031$   |
|    |                   | 128 (-)         | $2.779 {\pm} 0.015$ | $25.57 {\pm} 0.02$ | 36.08±0.98         | $0.8261 {\pm} 0.0019$ |
|    |                   | 64 (↓)          | 3.860±0.237         | $24.15 {\pm} 0.26$ | 45.13±2.79         | $0.7942{\pm}0.0154$   |
|    |                   | 32 (↓)          | $6.289 {\pm} 0.709$ | $22.04 \pm 0.47$   | $72.85{\pm}16.69$  | $0.7338 {\pm} 0.0225$ |

1133 Table 10: Reconstruction performance of data-driven RAQ-VAE (K = 128) with or without crossforcing on the CelebA test dataset

