

codebook collapse를 피하기 위한 stochastic quantization.

* Introduction.

VQ 방식에는 codebook collapse 문제가 대두되는데,

이를 해결하기 위해 EMA, codebook reset, tuning이 사용된다.

본문에서는 양자화의 문제를 보고, stochastic quantization을 살펴보자.

↳ 일반적인 VAE framework 내에서 훈련 가능

학습 중에 quantization process의 stochasticity \rightarrow self-annealing

\rightarrow CE의 성능 층적이 안정화됨 Gaussian 대신 vMF(von Mises-Fisher) 확률

Contribution.

1. SQ-VAE 제작

2. annealing으로 codebook usage ↑

3. 더욱 SQ-VAE 개선. (Gaussian, vMF)

4. 성능 향상 및 평가.

* Background.

VAE와 VQVAE는 모두 reconstruction loss와 latent penalty로 구성된다고 말할 수 있다.

* Stochastically quantized VAE

SQ-VAE는 VAE와 VQVAE의 혼합을 연결합니다.

↳ heuristic tech & hyperparam의 문제를 풀었습니다.

혼련 중인 정밀 quantization이 균형을 찾았습니다.

- Overview of SQ-VAE

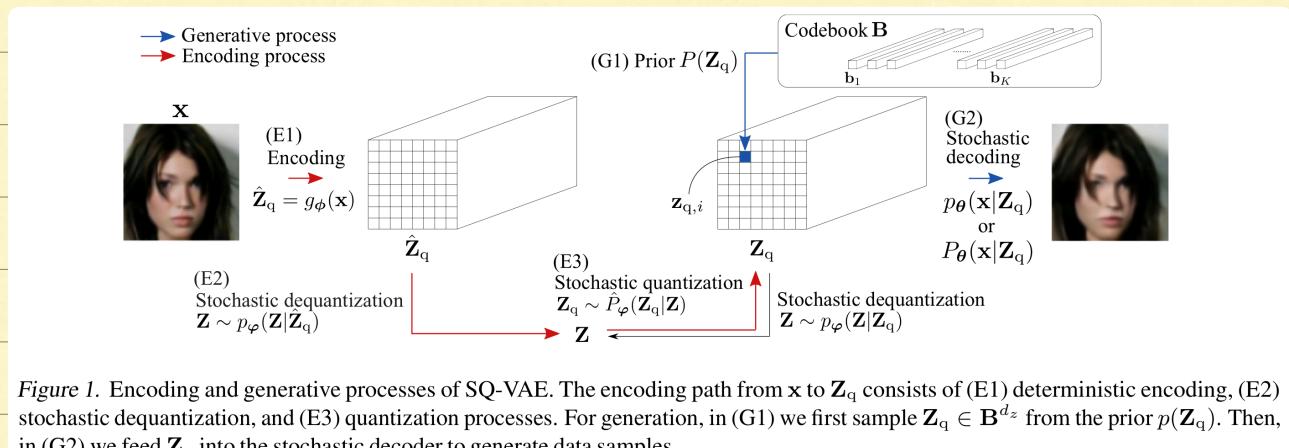


Figure 1. Encoding and generative processes of SQ-VAE. The encoding path from x to Z_q consists of (E1) deterministic encoding, (E2) stochastic dequantization, and (E3) quantization processes. For generation, in (G1) we first sample $Z_q \in \mathbb{B}^{d_z}$ from the prior $p(Z_q)$. Then, in (G2) we feed Z_q into the stochastic decoder to generate data samples.

Codebook을 B라고 할 때, $B := \{b_k\}_{k=1}^K$

그리고 G는 $x \sim P(x|Z_q)$, $x \sim P_{\text{data}}$, $Z_q \sim P(Z_q)$

$P(Z_q)$ 은 uniform i.i.d. (확률이 다 1/K)

Z_q 는 Z_q 로부터 dequantization 된 것이고,

$$P_\varphi(z|z_q), \text{ quantization: } \hat{P}_\varphi(z_q|z), z|z_q \sim P_\varphi(z|z_q)$$

latent
codebook depth

$$\hat{P}_\varphi(z_q|z) \propto P_\varphi(z|z_q) P(z_q)$$

$$\hat{z}_q = g_\phi(x), g_\phi: \mathbb{R}^D \rightarrow \mathbb{R}^{d_b \times d_z}$$

$\hat{z}_q \leftrightarrow z_q$ 가 가능합니다.

$P_\varphi(z|z_q)$ 과 $\hat{P}_\varphi(z_q|z)$ 을 쌓아, \hat{z}_q 와 z_q 를 연결

$Q_\omega(z_q|x) := E_{q_\omega(z|x)} [\hat{P}_\varphi(z_q|z)] ; \omega = \{\phi, \varphi\}, q_\omega(z|x) := P_\varphi(z|g_\phi(x))$ 이, x 부터 z_q 까지 stochastic encoding

SQ-VAE로 볼 때의 ELBO는

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &\geq -\mathcal{L}_{\text{SQ}}(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\omega}, \mathbf{B}) := \\ &= \mathbb{E}_{q_{\boldsymbol{\omega}}(\mathbf{Z}|\mathbf{x})\hat{P}_{\boldsymbol{\varphi}}(\mathbf{Z}_q|\mathbf{Z})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{Z}_q)p_{\boldsymbol{\varphi}}(\mathbf{Z}|\mathbf{Z}_q)\tilde{P}(\mathbf{Z}_q)}{q_{\boldsymbol{\omega}}(\mathbf{Z}|\mathbf{x})\hat{P}_{\boldsymbol{\varphi}}(\mathbf{Z}_q|\mathbf{Z})} \right] \\ &= \mathbb{E}_{q_{\boldsymbol{\omega}}(\mathbf{Z}|\mathbf{x})\hat{P}_{\boldsymbol{\varphi}}(\mathbf{Z}_q|\mathbf{Z})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{Z}_q)p_{\boldsymbol{\varphi}}(\mathbf{Z}|\mathbf{Z}_q)}{q_{\boldsymbol{\omega}}(\mathbf{Z}|\mathbf{x})} \right] \\ &\quad + \mathbb{E}_{q_{\boldsymbol{\omega}}(\mathbf{Z}|\mathbf{x})} H(\hat{P}_{\boldsymbol{\varphi}}(\mathbf{Z}_q|\mathbf{Z})) + \text{const.}, \end{aligned} \quad (5)$$

→ uniform 이므로 constant.

→ $H(P) \in \text{Pol entropy}$

↳ E.G. B) 품질의 optimizing

- Gaussian SQ-VAE

Gaussian SQ-VAE는 dequantization을 gaussian으로 가정.

$$p_{\boldsymbol{\varphi}}(\mathbf{z}_i|\mathbf{Z}_q) = \mathcal{N}(\mathbf{z}_{q,i}, \underline{\Sigma_{\boldsymbol{\varphi}}}), \quad (6)$$

trainable

Inverse of the quantization은 ($\mathbf{z}_q \Leftarrow \mathbf{z}$)

$$\hat{P}_{\boldsymbol{\varphi}}(\mathbf{z}_{q,i} = \mathbf{b}_k | \mathbf{z}) = \text{softmax}_k \left(\left[-\frac{(\mathbf{b}_j - \mathbf{z}_i)^T \Sigma_{\boldsymbol{\varphi}}^{-1} (\mathbf{b}_j - \mathbf{z}_i)}{2} \right]_{j=1}^k \right)$$

↳ Variance $\Sigma_{\boldsymbol{\varphi}}$ 와 \mathbf{z}_i 로부터 \mathbf{b}_j 의 Mahalanobis 거리를 unnormalized log-probability

Table 1. Different parameterizations of the variance $\Sigma_{\boldsymbol{\varphi}}$ in Gaussian SQ-VAE.

	Variance $\Sigma_{\boldsymbol{\varphi}}$	Unnormalized log-probability	Regularization objective $\mathcal{R}_{\boldsymbol{\varphi}}^{\mathcal{N}}(\mathbf{Z}, \mathbf{Z}_q)$
(I)	$\sigma_{\boldsymbol{\varphi}}^2 \mathbf{I}$	$\ \mathbf{b}_k - \mathbf{z}_i(\mathbf{x})\ _2^2 / 2\sigma_{\boldsymbol{\varphi}}^2$	$\ \mathbf{Z} - \mathbf{Z}_q\ _F^2 / 2\sigma_{\boldsymbol{\varphi}}^2$
(II)	$\sigma_{\boldsymbol{\varphi}}^2(\mathbf{x}) \mathbf{I}$	$\ \mathbf{b}_k - \mathbf{z}_i(\mathbf{x})\ _2^2 / 2\sigma_{\boldsymbol{\varphi}}^2(\mathbf{x})$	$\ \mathbf{Z} - \mathbf{Z}_q\ _F^2 / 2\sigma_{\boldsymbol{\varphi}}^2(\mathbf{x})$
(III)	$\sigma_{\boldsymbol{\varphi}, i}^2(\mathbf{x}) \mathbf{I}$	$\ \mathbf{b}_k - \mathbf{z}_i(\mathbf{x})\ _2^2 / 2\sigma_{\boldsymbol{\varphi}, i}^2(\mathbf{x})$	$\sum_{i=1}^{d_z} \ \mathbf{z}_i(\mathbf{x}) - \mathbf{z}_{q,i}\ _2^2 / 2\sigma_{\boldsymbol{\varphi}, i}^2(\mathbf{x})$
(IV)	$\text{diag}(\sigma_{\boldsymbol{\varphi}, i}^2(\mathbf{x}))$	$\sum_{j=1}^{d_b} (b_{k,j} - z_{i,j}(\mathbf{x}))^2 / 2\sigma_{\boldsymbol{\varphi}, i,j}^2(\mathbf{x})$	$\sum_{i=1}^{d_z} \sum_{j=1}^{d_b} (z_{i,j}(\mathbf{x}) - z_{q,i,j})^2 / 2\sigma_{\boldsymbol{\varphi}, i,j}^2(\mathbf{x})$

Decoding : $p_{\theta}(\mathbf{x}|\mathbf{z}_q) = \mathcal{N}(\mathbf{z}_{\theta}(\mathbf{z}_q), \sigma^2 \mathbf{I})$, $\sigma^2 \in \mathbb{R}_+$ 와 θ 는 trainable.

Encoding : $p_{\boldsymbol{\varphi}}(\mathbf{z}_i|\hat{\mathbf{Z}}_q) = \mathcal{N}(\hat{\mathbf{z}}_{q,i}, \Sigma_{\boldsymbol{\varphi}})$.

Objective function.

eq.5를 살펴보면.

$$\begin{aligned} \mathcal{L}_{N-\text{SE}} &= \mathbb{E}_{q_{\omega}(\mathbf{Z}|\mathbf{x}) \hat{P}_{\varphi}(\mathbf{Z}_q|\mathbf{Z})} \left[\frac{1}{2\sigma^2} \|\mathbf{x} - f_{\theta}(\mathbf{Z})\|_2^2 + \underline{\mathcal{R}_{\varphi}^N(\mathbf{Z}, \mathbf{Z}_q)} \right] \\ &\quad - \mathbb{E}_{q_{\omega}(\mathbf{Z}|\mathbf{x})} H \left(\hat{P}_{\varphi}(\mathbf{Z}_q|\mathbf{Z}) \right) + \frac{D}{2} \log \sigma^2 + \text{const.}, \quad (8) \end{aligned}$$

table 1의 regularization objective.

Appendix B.1

- Self-annealed Quantization.

Table 1의 type I을 살펴보자.

Σ_p 는 훈련 과정에서 quantization의 정도를 조절한다.

Appendix D.2의 $\sigma^2 \rightarrow \infty$ 및 $\sigma^2 \rightarrow 0$ 의 경우 extreme case를 살펴보자.

Proposition 1

$P_{\text{data}}(\mathbf{x})$ 가 finite하지만, \mathbf{g}_{θ} 와 $\{b_k\}_{k=1}^K$ 가 bounded라고 가정.

$w^* = \{\phi^*, \psi^*\}$ 가 fixed $\theta, \sigma^2, \{b_k\}_{k=1}^K$ 의 경우, $E_{P_{\text{data}}(\mathbf{x})} D_{\text{KL}}(Q_{\omega}(\mathbf{Z}_q|\mathbf{x}) \| P_{\theta}(\mathbf{Z}_q|\mathbf{x}))$ 의 minimizer라고 하면,

$\sigma^2 \rightarrow 0$ 이면, $\sigma_{\phi}^2 \rightarrow 0$

$\sigma^2 \rightarrow \infty$ 일 때, eq.8의 첫번째 term은 사라짐. $\rightarrow P_{\varphi}(\mathbf{Z}_{q,i} = b_k | \mathbf{Z})$ 가 uniform이 되기 때문.

$\sigma^2 \rightarrow 0$ 이면, kronecker delta $\delta_{k,i}$ 가 수렴 $\delta_{k,i}$. $\hat{k} = \arg \min_k \|\mathbf{z}_i - \mathbf{b}_k\|_2$

In quantization의 VAE의 posterior가 끝나.

$\rightarrow \sigma^2 \rightarrow \infty$ 이면, quantization process도 확률을 감소시키고, 예상치가 된다 \rightarrow self-annealing.

Dynamics of the variance parameter

self-annealing을 verify하기 위해. $\Sigma_p = \sigma_p^2 \mathbf{I}$ 의 gaussian SQ-VAE를 고려해보면,

σ_p^2 를 살펴보기 위해 σ_p^2 를 고정. \rightarrow Appendix D.3.

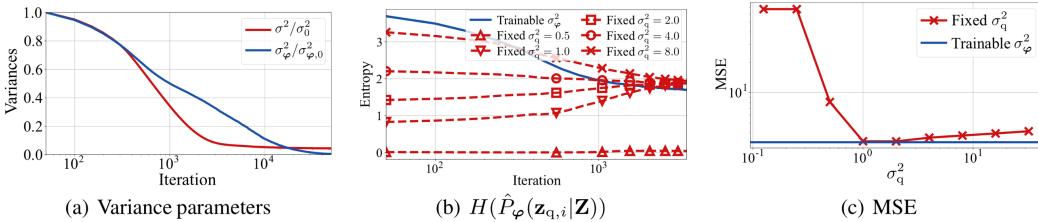


Figure 2. Empirical study on the dynamics related to σ_q^2 in Section 3.3. (a) The variance parameter σ_φ^2 (blue) decreased with σ^2 (red), where σ_0^2 and $\sigma_{\varphi,0}^2$ are their initial values. (b) Average entropy of the quantization process w.r.t. the iteration, which is obtained by Monte Carlo estimation. (c) MSE for trainable σ_φ^2 and various values of σ_q^2 on the test set.

- vMF SQ-VAE for categorical Distributions.

Categorical classification \Rightarrow softmax (softmax 사용)

CE or cross ELBO.

$$\begin{aligned} \mathcal{L}_{\text{CE-SQ}}^{\text{naive}} = & \mathbb{E}_{q_{\omega}(\mathbf{Z}|\mathbf{x}) \hat{P}_{\varphi}(\mathbf{Z}_q|\mathbf{Z})} \left[- \sum_{d=1}^D \log(P_{\theta}(x_d = c|\mathbf{Z}_q)) \right. \\ & \left. + \mathcal{R}_{\varphi}^{\mathcal{N}}(\mathbf{Z}, \mathbf{Z}_q) \right] - \mathbb{E}_{q_{\omega}(\mathbf{Z}|\mathbf{x})} H \left(\hat{P}_{\varphi}(\mathbf{Z}_q|\mathbf{Z}) \right) + \text{const.} \end{aligned} \quad (9a)$$

with

$$P_{\theta}(x_d = c|\mathbf{Z}_q) = \text{softmax}_c \left(\{\mathbf{w}_{\text{last},c'}^\top \tilde{f}_{\theta^-,d}^{\text{rest}}(\mathbf{Z}_q)\}_{c'=1}^{C_{\text{all}}} \right). \quad (9b)$$

$$\tilde{f}_{\theta^-,d}^{\text{rest}}(\mathbf{Z}_q) = \mathcal{B}^{d_x} \rightarrow \mathcal{R}^F$$

주의 경고, eq.8.24 다음에 그림에 대한 것과 같이 고려해야 한다. self-annealing의 효과를 넓힐 수 있다.

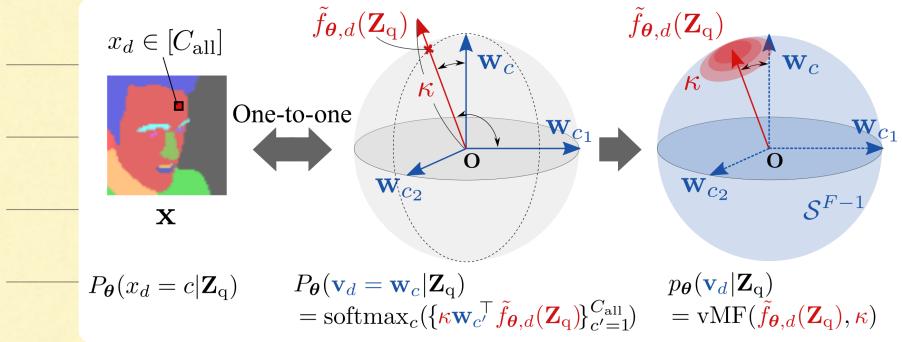
self-annealing을 위해, vMF 분포와 model을 refine.

F-dimension embedded 된 hyper-sphere S^{F-1} 을 사용.

$W_c \in S^{F-1}$ 의 surface with, C-dim data categorical의 projected vector를 명시연.

projected data $x_{d,c}$ hypersphere에 속하는 $V_d \in \{W_c\}_{c=1}^{C_{\text{all}}}$

$x_d = c$ 인 경우, $V_d = W_c$, $V_d = W_c$ 면 $x_d = c$



Decoding

$$P_{\theta}(\mathbf{v}_d = \mathbf{w}_c | Z_q) = \text{softmax}_c \left(\left\{ \kappa \mathbf{w}_{c'}^\top \tilde{f}_{\theta,d}(Z_q) \right\}_{c'=1}^{C_{\text{all}}} \right), \quad (10)$$

$$p_{\theta}(\mathbf{v}_d | Z_q) \propto \exp \left(\kappa \mathbf{v}_d^\top \tilde{f}_{\theta,d}(Z_q) \right). \quad (11)$$

$$p_{\theta}(\mathbf{v}_d | Z_q) = \text{vMF}(\tilde{f}_{\theta,d}(Z_q), \kappa)$$

Encoding

vMF를 사용하여, dequantization을 modeling.

$$p_{\varphi}(\mathbf{z}_i | Z_q) = \text{vMF}(\mathbf{z}_{q,i}, \kappa_{\varphi}), \quad (12)$$

κ_{φ} trainable concentration parameter.

$$\hat{P}_{\varphi}(\mathbf{z}_{q,i} = \mathbf{b}_k | Z) = \text{softmax}_k \left(\{\kappa_{\varphi} \mathbf{b}_j^\top \mathbf{z}_i\}_{j=1}^K \right), \quad (13)$$

Objective function.

$$\begin{aligned} \mathcal{L}_{\text{vMF-SQ}} &= \mathbb{E}_{q_{\omega}(\mathbf{Z} | \mathbf{x}) \hat{P}_{\varphi}(\mathbf{Z}_q | \mathbf{Z})} \left[-\kappa \sum_{d=1}^D \mathbf{v}_d^\top \tilde{f}_{\theta,d}(\mathbf{Z}_q) + \mathcal{R}_{\varphi}^{\text{vMF}}(\mathbf{Z}, \mathbf{Z}_q) \right] \\ &\quad - \mathbb{E}_{q_{\omega}(\mathbf{Z} | \mathbf{V})} H \left(\hat{P}_{\varphi}(\mathbf{Z}_q | \mathbf{Z}) \right) - \log C_F(\kappa) + \text{const.}, \end{aligned} \quad (14)$$

$$\mathcal{R}_{\varphi}^{\text{vMF}}(\mathbf{Z}, \mathbf{Z}_q) = \sum_{i=1}^{d_z} \kappa_{\varphi,i} (1 - \mathbf{z}_{q,i}^\top \mathbf{z}_i)$$

Comparing vMF SQ-VAE with Naïve Categorical SQ-VAE

SQ-VAE ≈ self-annealing SVAE

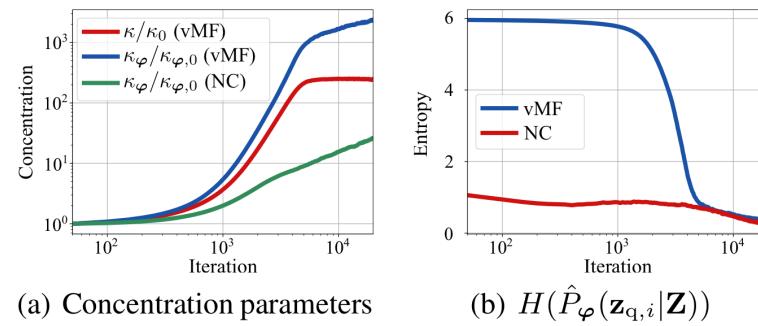


Figure 4. Comparison between vMF and NC decoders: (a) The concentration parameter of vMF decoder κ_φ increases with κ , whereas the growth of κ_φ of the NC decoder is relatively small. Here, κ_0 and $\kappa_{\varphi,0}$ indicate initial values. (b) Average entropy of quantization processes.

* Experiments

Continuous data Distribution.

Vision

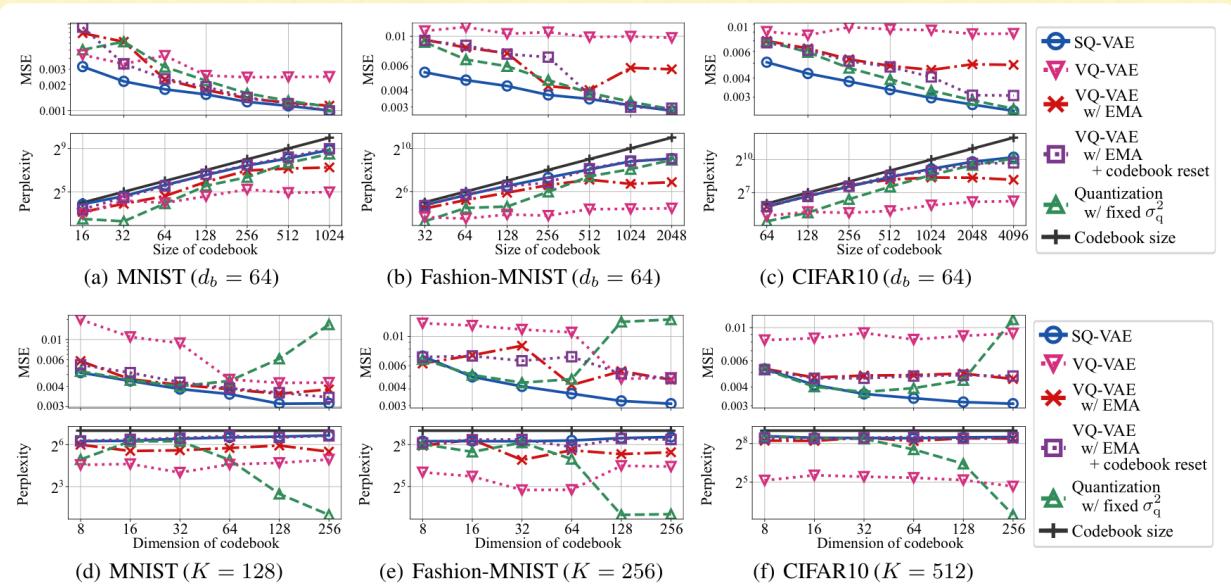


Figure 5. Empirical studies on the impact of codebook capacity examined on MNIST Fashion-MNIST and CIFAR10. (a)-(c) The size K is swept with the dimension d_b fixed to 64. (d)-(f) Various d_b values are tested with the size K fixed as 128, 256, and 512, respectively. The black lines with "+" marks indicate the upper bounds of the perplexities, i.e., K . All the y-axes are in log-scale.

Table 2. Evaluation on CelebA. The MSE ($\times 10^3$) and reconstructed FID (rFID) are evaluated using the test set. The codebook capacity for the discrete latent space is set to $(n_b, k) = (64, 512)$. The Roman numerals for Gaussian SQ-VAEs correspond to those in Table 1. We also show the FID of samples generated with a prior learned with PixelCNN.

Model	Reconstruction		Generation (FID)	Latent manipulation (FID)			
	MSE	rFID		Neighbor-3	Neighbor-5	Neighbor-10	Interp.
VAE	4.79 \pm 0.01	40.3 \pm 0.3	—	—	—	—	—
VQ-VAE (EMA)	1.33 \pm 0.41	18.5 \pm 5.1	42.0 \pm 11.5	31.9 \pm 14.8	42.8 \pm 20.7	70.7 \pm 35.4	28.2 \pm 6.4
VQ-VAE (EMA+code reset)	1.62 \pm 0.36	22.0 \pm 5.9	51.8 \pm 10.8	39.7 \pm 12.0	52.7 \pm 14.7	83.2 \pm 20.4	32.6 \pm 7.1
Quantization w/ fixed σ_q^2	1.09 \pm 0.01	15.9 \pm 0.1	38.2 \pm 0.9	20.0 \pm 0.4	26.4 \pm 0.8	41.5 \pm 2.1	18.6 \pm 0.3
Gaussian SQ-VAE (I)	0.96 \pm 0.01	14.8 \pm 0.3	28.2 \pm 0.9	17.8 \pm 0.1	21.9 \pm 0.1	33.1 \pm 0.3	17.6 \pm 0.6
Gaussian SQ-VAE (II)	0.98 \pm 0.01	14.3 \pm 0.2	27.7 \pm 1.1	17.8 \pm 0.2	22.2 \pm 0.4	34.0 \pm 0.9	17.6 \pm 0.1
Gaussian SQ-VAE (III)	0.96 \pm 0.00	13.9 \pm 0.1	28.1 \pm 0.3	17.3 \pm 0.2	21.6 \pm 0.3	33.5 \pm 0.6	18.5 \pm 0.4

Speech .

Table 3. Evaluation on VCTK and ZeroSpeech 2019. The MSE (dB²) of sample reconstruction is evaluated using the test set. We do not apply SQ-VAE (II) in this evaluation because of the variable length property of speech data and the different manipulations of speech signals between training and inference (see Appendix E.2).

Model	MSE (dB ²)	
	VCTK	ZeroSpeech 2019
VQ-VAE w/ EMA	29.59 \pm 0.25	34.33 \pm 1.57
Gaussian SQ-VAE (I)	25.52 \pm 0.08	33.17 \pm 1.11
Gaussian SQ-VAE (III)	25.94 \pm 0.22	34.35 \pm 1.07
Gaussian SQ-VAE (IV)	24.68 \pm 0.21	32.32 \pm 0.88

Categorical Distribution.

Table 4. Evaluation on CelebA-Mask. The pixel error (%), mIoU, and perplexity are evaluated using the test set.

Model	Pixel error	mIoU	Perplexity
VAE	8.79 \pm 0.01	55.8 \pm 0.3	—
VQ-VAE w/ EMA	6.95 \pm 0.14	59.7 \pm 0.7	46.2 \pm 2.0
NC SQ-VAE	6.63 \pm 1.38	64.1 \pm 5.4	12.6 \pm 5.2
vMF SQ-VAE	3.51 \pm 0.17	74.6 \pm 0.0	52.4 \pm 0.8

Table 5. Evaluation on MNIST and gray-scaled CelebA. The MSE ($\times 10^3$) is evaluated using the test set.

Model	MNIST	Gray-CelebA
VAE	22.80 \pm 0.32	17.73 \pm 0.23
VQ-VAE w/ EMA	6.24 \pm 0.18	5.19 \pm 0.06
NC SQ-VAE	10.89 \pm 0.47	3.88 \pm 0.02
vMF SQ-VAE	1.63 \pm 0.21	2.37 \pm 0.01

➤ Conclusion.

VAE의 VAE를 통해 dequantized와 quantized로 쌍으로 학습 CL.