

긴 AR은 long-dependency가 약하고, 시퀀스도 오래 걸린다.

∴ RQ-VAE는 VQ stack

RQ-Transformer는 다음 step을 예측하는 방식,

### \* Introduction.

VQ로 양자화한 후, AR로 sequence 생성.

AR 길이를 줄이는 것이 시간 계산이 속도화되지만, 성능과 trade-off

그리고 너무 큰 codebook은 collapse 유발.

RQ로 feature map을 더 잘 균사.

↳ Coarse 뒤에 fine까지 재구성으로.

RQ-Transformer는 양자화된 feature map or feature vector를 hidden input.

↳ DNN의 코드를 한번에 계산.

↳ soft labeling, stochastic sampling 추가.

### Contribution.

1) stack VQ가 가능한 RQ-VAE 개발.

2) RQ-Transformer로, exposure bias 없음, codebook 크기 ↓

## Method.

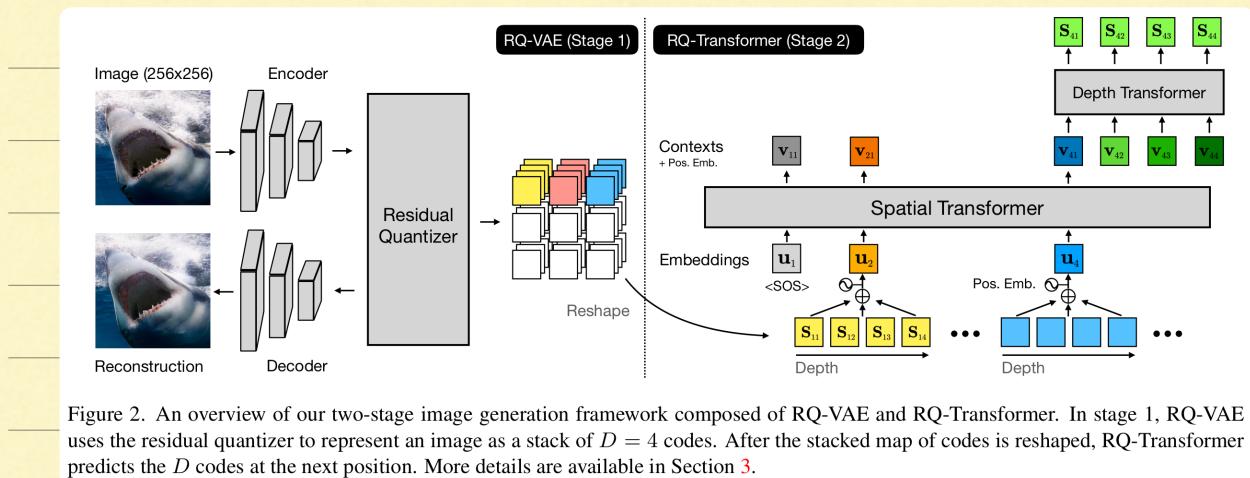


Figure 2. An overview of our two-stage image generation framework composed of RQ-VAE and RQ-Transformer. In stage 1, RQ-VAE uses the residual quantizer to represent an image as a stack of  $D = 4$  codes. After the stacked map of codes is reshaped, RQ-Transformer predicts the  $D$  codes at the next position. More details are available in Section 3.

## \* Stage 1

VQVAE 공식 소개 후, codebook의 크기를 늘리지 않고도 feature map을 더 정확히 고사하는 RQ-VAE 소개.

### - Formulation.

Codebook  $C$ 의 경우  $\{k, e(k)\}_{k \in [K]}$ ,  $e(k) \in \mathbb{R}^{n_z}$ ,  $n_z$ 는 codebook dim.,  $k$ 는 index.

$Q(z; C)$ 는 \$z\$가 Quantize 된 것을 의미.

$$Q(z; C) = \arg \min_{k \in [K]} \|z - e(k)\|_2^2. \quad (1)$$

\$\Rightarrow\$ feature vector이 \$Q\$ 적용. \$\rightarrow\$ code map \$M \in K^{H \times W}\$, quantized feature map \$\hat{z} \in \mathbb{R}^{H \times W \times n\_z}\$

$$M_{hw} = Q(Z_{hw}; C), \quad \hat{z}_{hw} = e(M_{hw}), \quad (2)$$

$$\text{recon: } \hat{x} = G(\hat{z})$$

$M$ 이  $H, W$ 인 ARI cost per time의 중요.

그러나 VQVAE는 lossy compression이 가능하다. bit =  $HW \log_2 k$ .

\$\therefore H/L, W/L\$인 recon loss를 유지하면서면,  $k^4$ 가 필요.

↳ 하지만  $k$ 를 커으면 codebook collapse가 일어남.

- Residual Quantization.

코드북 크기를 늘리는 대신 RQ. depth: D.

$$RQ(\mathbf{z}; \mathcal{C}, D) = (k_1, \dots, k_D) \in [K]^D, \quad (3)$$

RQ는 차례로  $k_d$  만큼 recursively

$$\begin{aligned} k_d &= Q(\mathbf{r}_{d-1}; \mathcal{C}), \rightarrow \text{index} \\ \mathbf{r}_d &= \mathbf{r}_{d-1} - \mathbf{e}(k_d), \rightarrow \text{시작된 다음 r_d}. \end{aligned} \quad (4)$$

Vector를 coarse 를 fine 까지.

$\hat{\mathbf{z}}^{(c)}$ 이  $e(k_i)$  과 가장 가깝고, feature에도 가까울

나머지는 error를 줄이는 데 사용 → 부분적인 더듬이 더 나은 근사제공.

더 많은 codebook은 더 공유한. → 더 커지 이득.

1. 코드북 크기가 작아짐.

2. utility↑.

즉, 더 정확한 근사

VQVAE의 경우  $R^m$  dim으로 k-clustering이지만 RQ는  $K^D$  clustering

↳ D를 늘리면 exponentially increase.

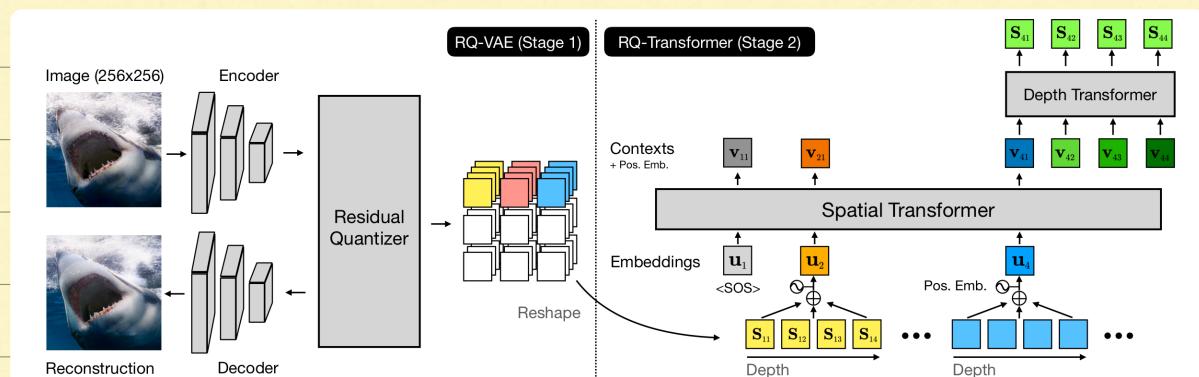


Figure 2. An overview of our two-stage image generation framework composed of RQ-VAE and RQ-Transformer. In stage 1, RQ-VAE uses the residual quantizer to represent an image as a stack of  $D = 4$  codes. After the stacked map of codes is reshaped, RQ-Transformer predicts the  $D$  codes at the next position. More details are available in Section 3.

## - RQ-VAE

VQVAE  $\rightarrow$  RQ

code M:  $[k]^{H \times W \times D}$ ,  $\hat{\mathbf{Z}}^{(d)} \in \mathbb{R}^{H \times W \times n_d}$

$$\mathbf{M}_{hw} = \mathcal{RQ}(E(\mathbf{X})_{hw}; \mathcal{C}, D),$$

$$\hat{\mathbf{Z}}_{hw}^{(d)} = \sum_{d'=1}^d \mathbf{e}(\mathbf{M}_{hwd'}). \quad (5)$$

down sample factor f는 는 2(2, (H,W)는 높이와 너비,

AR의 계산은 높이를 기준으로 한다.

## - Training RQ-VAE.

$$\mathcal{L}_{\text{recon}} = \|\mathbf{X} - \hat{\mathbf{X}}\|_2^2, \quad (6)$$

$$\mathcal{L}_{\text{commit}} = \sum_{d=1}^D \left\| \mathbf{Z} - \text{sg} \left[ \hat{\mathbf{Z}}^{(d)} \right] \right\|_2^2, \quad (7)$$

$$\mathcal{L} = \mathcal{L}_{\text{recon}} + \beta \mathcal{L}_{\text{commit}}$$

Codebook은 coarse-to-fine으로 feature map 구조.

## \* Stage 2. RQ - Transformer.

RQ codebook을 학습하는 과정

- AR modeling for Codes for depth D.

AR의 순서를 raster로 볼 때 Sequence의 길이  $S \in [k]^{T \times D}$ ,  $T = HW$

$$\mathbf{S}_t = (\mathbf{S}_{t1}, \dots, \mathbf{S}_{tD}) \in [K]^D \quad \text{for } t \in [T]. \quad (8)$$

AR model은  $p(S)$ 를 학습

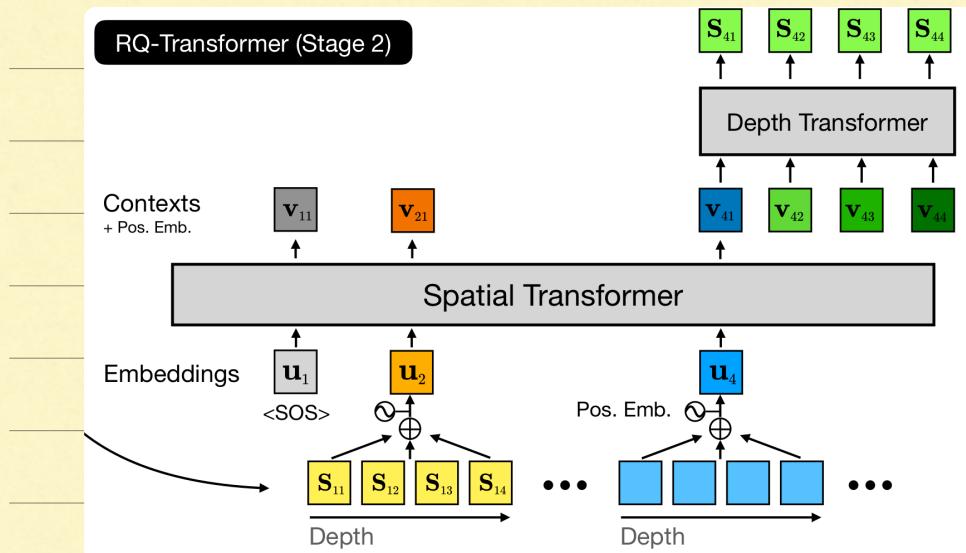
$$p(\mathbf{S}) = \prod_{t=1}^T \prod_{d=1}^D p(\mathbf{S}_{td} | \mathbf{S}_{<t,d}, \mathbf{S}_{t,<d}). \quad (9)$$

- RQ-Transformer Architecture.

TD개를 고려 생성하는 것은 length로 짚고 compute로 가볍고 빠르다.

$\therefore D$ 를 한번에 생성

In depth transformer, spatial transformer.



- Spatial transformer.

input의 경도 이전 embedding과 raster-scan number

$$\mathbf{u}_t = \text{PE}_T(t) + \sum_{d=1}^D \mathbf{e}(\mathbf{S}_{t-1,d}) \quad \text{for } t > 1, \quad (10)$$

$$\mathbf{h}_t = \text{SpatialTransformer}(\mathbf{u}_1, \dots, \mathbf{u}_t). \quad (11)$$

- Depth Transformer.

context vector  $h_t$ 을 계산하는 (t의 경우  $\mathbf{h}_t$ )

↑ t의 경우 depth를 계산

$$\mathbf{v}_{td} = \text{PE}_D(d) + \sum_{d'=1}^{d-1} \mathbf{e}(\mathbf{S}_{td'}) \quad \text{for } d > 1, \quad (12)$$

t의 경우 depth embedding을 계산하는 가정.

$$\therefore P_{td}(k) = p(S_{td=k} | S_{<t,d}, S_{t,<d}) \neq 0,$$

$$\mathbf{p}_{td} = \text{DepthTransformer}(\mathbf{v}_{t1}, \dots, \mathbf{v}_{td}). \quad (13)$$

↑ AR NLL loss는.

$$\mathcal{L}_{AR} = \mathbb{E}_{\mathbf{S}} \mathbb{E}_{t,d} [-\log p(\mathbf{S}_{td} | \mathbf{S}_{<t,d}, \mathbf{S}_{t,<d})]. \quad (14)$$

- Computational Complexity

TXD sequence를 만드는 데의 복잡성 complexity

$$O(NT^2D^2) \rightarrow O(N_r T^2 + N_d TD^2)$$

### - Soft Labeling and Stochastic Labeling

RQ는 오류가 누적되는 데 따른 exposure bias↑

soft labeling & stochastic sampling은 이를 해결하기 위한

$$Q_\tau(k|\mathbf{z}) \propto e^{-\|\mathbf{z} - \mathbf{e}(k)\|_2^2/\tau} \quad \text{for } k \in [K]. \quad (15)$$

### - Soft labeling of target code.

코드는 embedding 히든 기반

### - Stochastic sampling for codes of RQ-VAE

## \* Experiments.

Table 1. Comparison of FIDs for unconditional image generation on LSUN-{Cat, Bedroom, Church} [48] and FFHQ [25].

	Cat	Bedroom	Church	FFHQ
VDVAE [7]	-	-	-	28.5
DDPM [21]	19.75	4.90	7.89	-
ImageBART [13]	15.09	5.51	7.32	9.57
StyleGAN2 [26]	7.25	2.35	3.86	3.8
BigGAN [3]	-	-	-	12.4
DCT [33]	-	6.40	7.56	13.06
VQ-GAN [14]	17.31	6.35	7.81	11.4
<b>RQ-Transformer</b>	8.64	3.04	7.45	10.38

Table 3. Comparison of FID and CLIP score [36] on the validation data of CC-3M [43] for text-conditioned image generation.

	Params	FID	CLIP-s
VQ-GAN [14]	600M	28.86	0.20
ImageBART [13]	2.8B	22.61	0.23
<b>RQ-Transformer</b>	654M	12.33	0.26

Table 2. Comparison of FIDs and ISs for class-conditioned image generation on ImageNet [9] 256×256. † denotes a model without our stochastic sampling and soft labeling. ‡ denotes the use of rejection sampling or gradient guidance by pretrained classifier. \* denotes the use of RQ-VAE trained for 50 epochs.

	Params	FID	IS
<b>without rejection sampling or gradient guidance</b>			
ADM [11]	554M	10.94	101.0
ImageBART [13]	3.5B	21.19	61.6
BigGAN [3]	164M	7.53	168.6
BigGAN-deep [3]	112M	6.84	203.6
VQ-VAE2 [39]	13.5B	~31	~45
DCT [33]	738M	36.5	n/a
VQ-GAN [14]	1.4B	15.78	74.3
<b>RQ-Transformer</b>	480M	15.72	86.8±1.4
<b>RQ-Transformer</b> †	821M	14.06	95.8±2.1
<b>RQ-Transformer</b>	821M	13.11	104.3±1.5
<b>RQ-Transformer</b>	1.4B	11.56	112.4±1.1
<b>RQ-Transformer</b> *	1.4B	8.71	119.0±2.5
<b>RQ-Transformer</b> *	3.8B	7.55	134.0±3.0
<b>with rejection sampling or gradient guidance</b>			
ADM‡ [11]	554M	4.59	186.7
ImageBART‡ [13]	3.5B	7.44	273.5±4.1
VQ-VAE2‡ [39]	13.5B	~10	~330
VQ-GAN‡ [14]	1.4B	5.20	280.3±5.5
<b>RQ-Transformer</b> ‡	1.4B	4.45	326.0±3.5
<b>RQ-Transformer</b> *‡	1.4B	3.89	337.5±4.6
<b>RQ-Transformer</b> *‡	3.8B	3.80	323.7±2.8
Validation Data	-	1.62	234.0

Table 4. Comparison of FIDs between ImageNet validation images and their reconstructed images according to codebook size ( $K$ ) and the shape of code map  $H \times W \times D$ . † denotes the reproduced performance, and \* denotes 50 epochs of training.

	$H \times W \times D$	$K$	rFID
VQ-GAN [14]	$16 \times 16 \times 1$	16,384	4.90
VQ-GAN <sup>†</sup>	$16 \times 16 \times 1$	16,384	4.32
VQ-GAN	$8 \times 8 \times 1$	16,384	17.95
VQ-GAN	$8 \times 8 \times 1$	65,536	17.66
VQ-GAN	$8 \times 8 \times 1$	131,072	17.09
RQ-VAE	$8 \times 8 \times 2$	16,384	10.77
RQ-VAE	$8 \times 8 \times 4$	16,384	4.73
RQ-VAE*	$8 \times 8 \times 4$	16,384	3.20
RQ-VAE	$8 \times 8 \times 8$	16,384	2.69
RQ-VAE	$8 \times 8 \times 16$	16,384	1.83

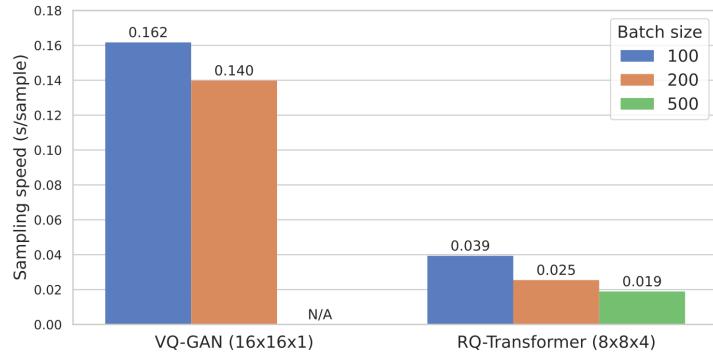


Figure 4. The sampling speed of RQ-Transformer with 1.4B parameters according to batch size and code map shape.

## → Ablation



Figure 5. The examples of coarse-to-fine approximation by RQ-VAE. The first example is the original image, and the others are reconstructed from  $\hat{\mathbf{Z}}^{(d)}$ . As  $d$  increases, the reconstructed images become clear and include fine-grained details of the original image.

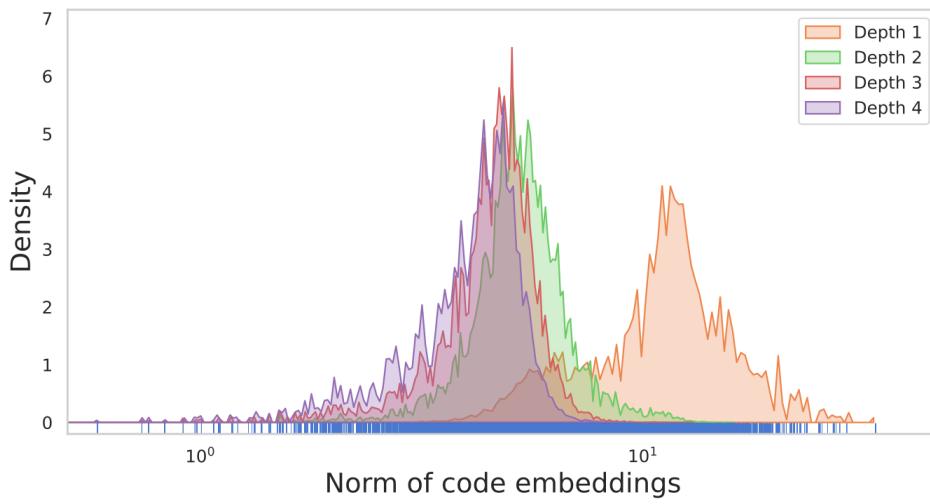


Figure 6. The distribution of used codes at each quantization depth. The blue bar plot represents the code distribution according to the norm of embeddings. ImageNet validation data is used.

## \* Appendix