# **Hierarchical Residual Learning Based Vector Quantized Variational** Autoencoder for Image Reconstruction and Generation

Mohammad Adiban<sup>1,2</sup>, Kalin Stefanov<sup>2</sup>, Marco Siniscalchi<sup>1</sup>, Giampiero Salvi<sup>1,3</sup>





- HR-VQVAE maps the continuous latent representations to several layers of ulletdiscrete representations through hierarchical codebooks.
- The vector selected within bottom layer determines the codebook that is  $\bullet$ activated in the top layer.
- Such a hierarchical searching procedure provides the advantage of local access lacksquareto codebook indexes, which dramatically reduces search time.

### **Objective functions**

Quantize<sup>*i*</sup>(
$$\boldsymbol{\xi}_{hw}^{i-1}$$
) =  $\mathbf{e}_{k}^{i}$  where  $k = \arg\min_{j} \|\boldsymbol{\xi}_{hw}^{i-1} - \mathbf{e}_{j}^{i}\|_{2}$ , (3)

$$\mathbf{e}_C = \sum_{i=1}^n \mathbf{e}^i,\tag{4}$$

KTH VETENSKAP

MONASH University

**BMVC** 

$$\mathcal{L}(\mathbf{x}, \mathcal{D}(\mathbf{e}_C)) = \|\mathbf{x} - \mathcal{D}(\mathbf{e}_C)\|_2^2 + \|\mathbf{sg}[\boldsymbol{\xi}^0] - \mathbf{e}_C\|_2^2 + \beta_0 \|\mathbf{sg}[\mathbf{e}_C] - \boldsymbol{\xi}^0\|_2^2 + \sum_{i=1}^n \mathcal{L}(\boldsymbol{\xi}^{i-1}, \mathbf{e}^i), \quad (5)$$

$$\mathcal{L}(\boldsymbol{\xi}^{i-1}, \mathbf{e}^{\mathbf{i}}) = \|\mathbf{sg}[\boldsymbol{\xi}^{i-1}] - \mathbf{e}^{i}\|_{2}^{2} + \beta_{i}\|\mathbf{sg}[\mathbf{e}^{i}] - \boldsymbol{\xi}^{i-1}\|_{2}^{2}, \tag{6}$$

## Experiments



- A novel objective function is proposed to provide contrastive learning by • pushing each layer to extract information not learned by its preceding layers.
- The objective optimizes the output image from the combination of  $\bullet$ representations obtained from all layers.

## Background: VQVAE, VQVAE-2

Overview



Quantize
$$(\mathbf{z}_{hw}) = \mathbf{e}_k$$
 where  $k = \arg\min_j ||\mathbf{z}_{hw} - \mathbf{e}_j||_2,$  (1)

$$\mathcal{L}(\mathbf{x}, \mathcal{D}(\mathbf{e})) = \|\mathbf{x} - \mathcal{D}(\mathbf{e})\|_2^2 + \|\mathbf{sg}[\mathbf{z}] - \mathbf{e}\|_2^2 + \beta \|\mathbf{sg}[\mathbf{e}] - \mathbf{z}\|_2^2.$$
(2)





Fig. 2. VQVAE-2 Architecture [2].

#### **Proposed Method: HR-VQVAE**



1110401	FFHQ	ImageNet	CIFAR10	MNIST
VQVAE [19]	2.86/0.00298	3.66/0.00055	21.65/0.00092	7.9/0.00041
VQVAE-2 [18]	1.92/0.00195	2.94/0.00039	<b>18.03</b> /0.00068	6.7/0.00025
HR-VQVAE	1.26/0.00163	2.28/0.00027	18.11/ <b>0.00041</b>	6.1/0.00011

Table 2: Time for reconstructing 10,000 samples using HR-VQVAE, VQVAE-2 and VQVAE.

Model	Seconds				
	FFHQ	Imagenet	CIFAR10	MNIST	
VQVAE [19]	5.0977652	4.6152677	2.7087896	0.062474	
VQVAE-2 [18]	9.3443758	8.8135872	4.4492340	0.090778	
HR-VQVAE	0.8398101	0.6714823	0.4667842	0.010830	



Fig. 8. Random samples on FFHQ generated by VQVAE-2, VQ-GAN and HR-VQVAE, respectively.

Table 3: Generation results using HR-VQVAE, VQVAE-2 and VQVAE.

Model	Generation evaluation (FID $\downarrow$ )			
	FFHQ	ImageNet	CIFAR10	MNIST
VQVAE [19]	24.93	44.76	78.90	16.69
VQVAE-2 [18]	19.66	39.51	74.43	11.81
HR-VQVAE	17.45	35.29	71.38	11.75

#### Conclusion

- We proposed a novel multi-layer variational Autoencoder method for image modeling that we call HR-VQVAE.
- HR-VQVAE learns hierarchical residual discrete representations in an iterative and hierarchical fashion.
- The objective of HR-VQVAE is designed to encourage different layers to encode different aspects of an image.
- Through experimental evidence, we show how HR-VQVAE can reconstruct images with a higher level of details than state-of-the-art models with similar complexity.
- We also show that we can increase the size of the codebooks without incurring the codebook collapse problem that is observed in methods such as VQVAE and VQVAE-2.

#### References

[1] Aaron van den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In Advances in Neural Information Processing Systems, pages 6306–6315, 2017.

[2] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. In Advances in neural information processing systems, pages 14866–14876, 2019.