

# Neural Discrete Representation Learning

Van den Oord et al.

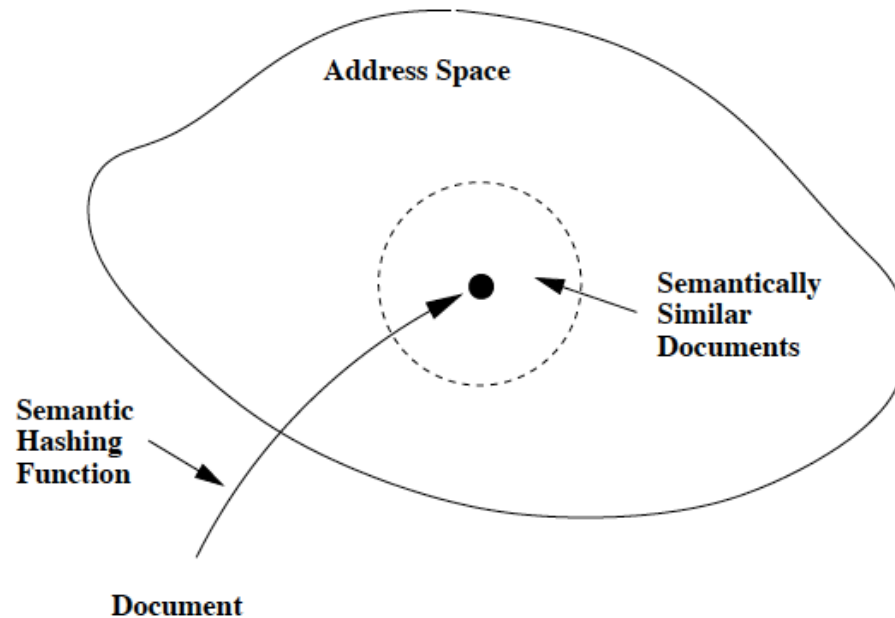
---

Park Jungsoo

Data Mining & Information Systems Lab.  
Department of Computer Science and Engineering,  
College of Informatics, Korea University

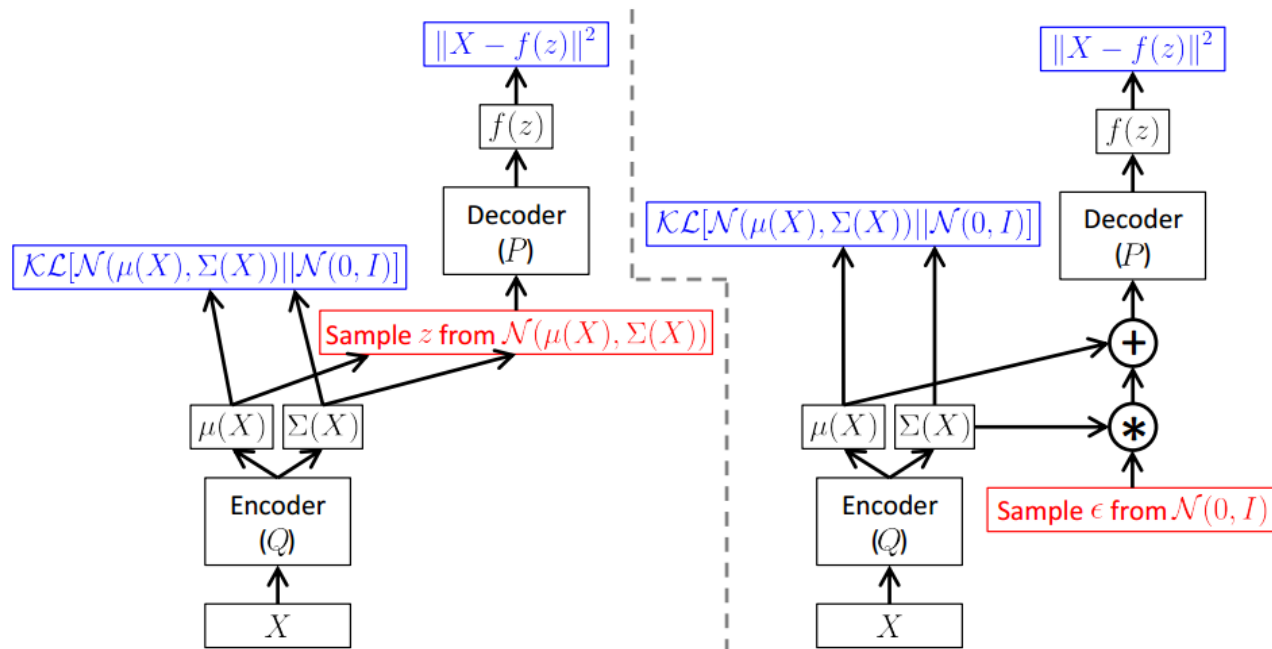
---

## Why Discrete Latent Representation?



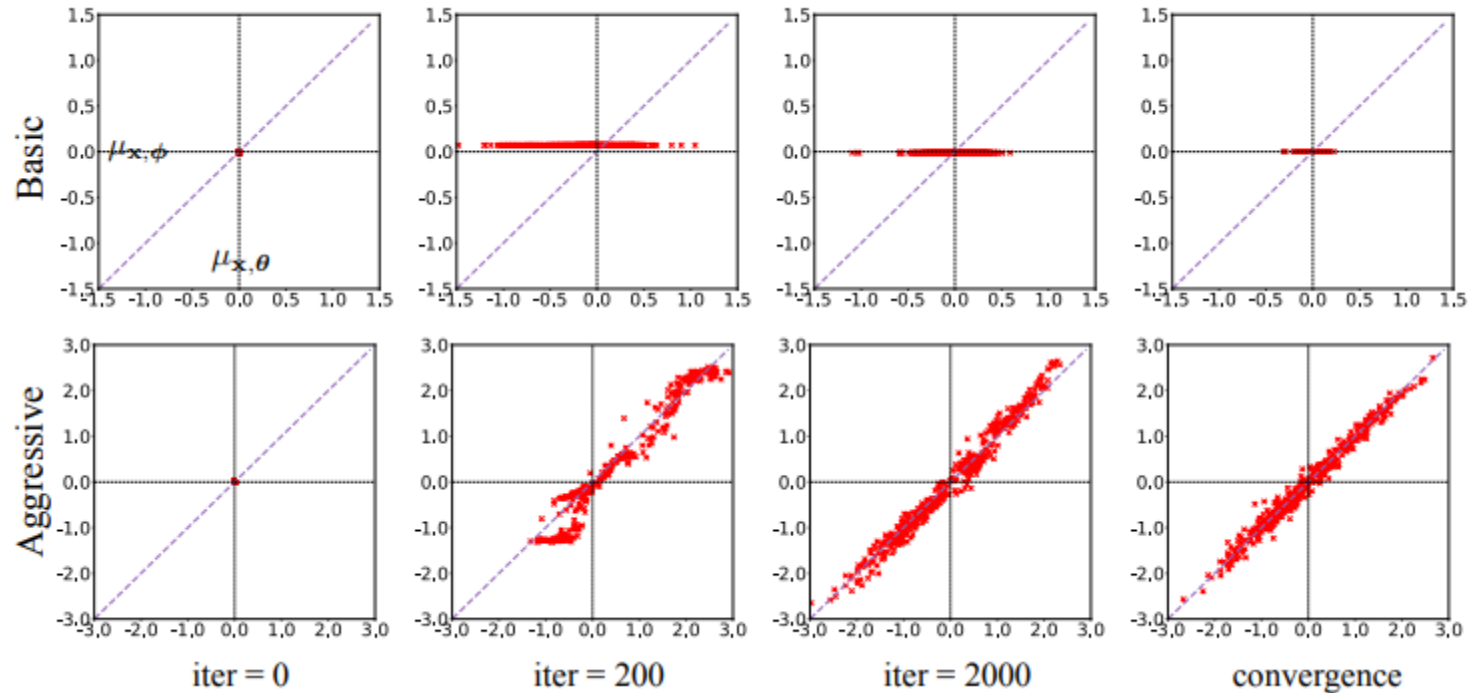
- Computational Efficiency
- Interpretability and Communication
- More Natural

## Auto-Encoding Variational Bayes



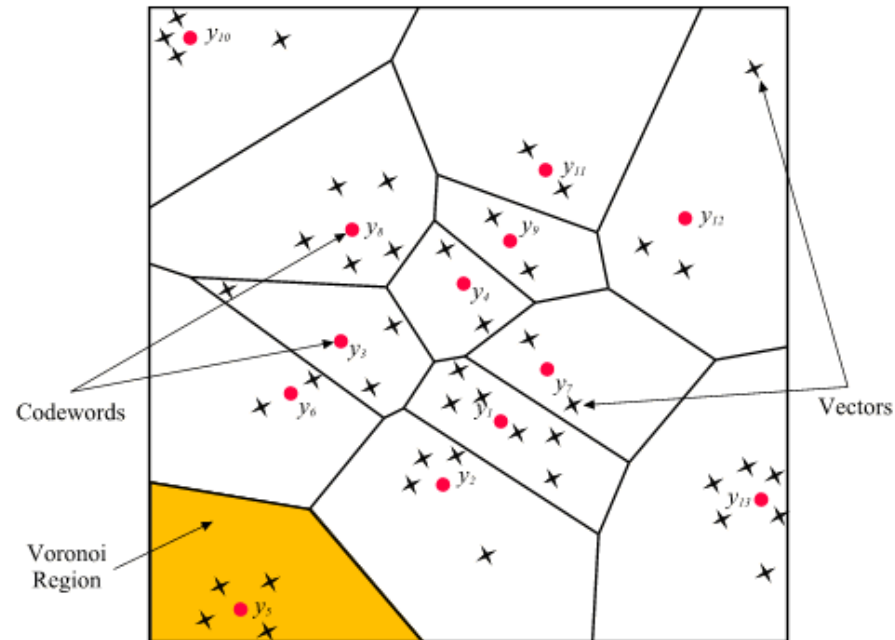
- Maximizing ELBO
- Ancestral sampling from Standard Normal D.

## Posterior Collapse



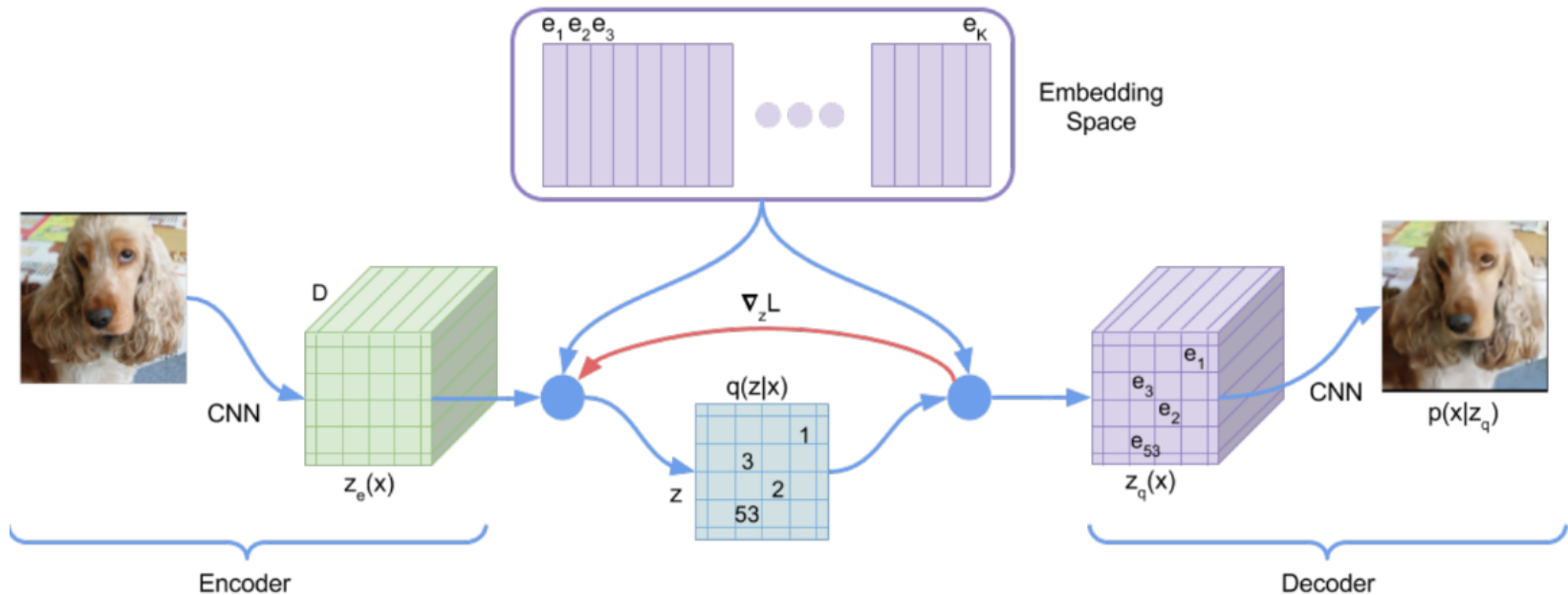
- With strong autoregressive decoder, posterior collapse happens
- Not capturing meaningful representation

## Vector Quantization



- Quantization technique for modeling probability density function by distribution of prototype vectors
- Originally used in lossy data compression

## Overview



## Categorical Distribution

- Posterior and Prior are categorical distributions.

$$L = \log p(x|z_q(x)) + \|\text{sg}[z_e(x)] - e\|_2^2 + \beta \|z_e(x) - \text{sg}[e]\|_2^2,$$

$$q(z = k|x) = \begin{cases} 1 & \text{for } k = \text{argmin}_j \|z_e(x) - e_j\|_2, \\ 0 & \text{otherwise} \end{cases},$$

$$z_q(x) = e_k, \quad \text{where } k = \text{argmin}_j \|z_e(x) - e_j\|_2$$

## Resolving Posterior Collapse

- As for vanilla VAE,

$$\begin{aligned} D_{KL}((q_{\phi}(\mathbf{z})||p_{\theta}(\mathbf{z}))) &= \int q_{\theta}(\mathbf{z}) (\log p_{\theta}(\mathbf{z}) - \log q_{\theta}(\mathbf{z})) d\mathbf{z} \\ &= \frac{1}{2} \sum_{j=1}^J (1 + \log((\sigma_j)^2) - (\mu_j)^2 - (\sigma_j)^2) \end{aligned}$$

- As for VQ VAE,

$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right).$$

For specific  $k$ ,  $p(x)$ , posterior is 1 and 0 otherwise, therefore yielding

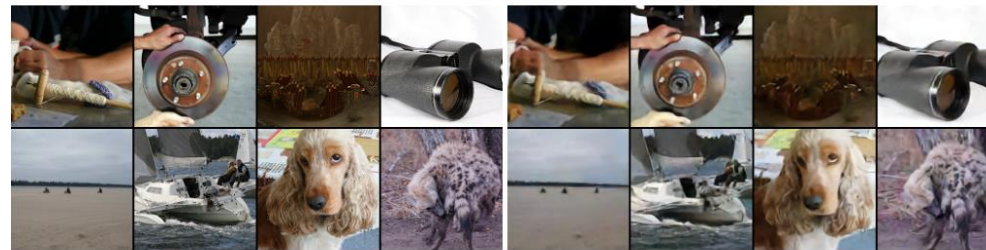
$$\log K$$

**Thus, posterior collapse(KL term  $\rightarrow 0$ ) problem doesn't happen**



## Autoregressive Prior

- After training with uniform prior, autoregressive model is fit for learning prior distribution, thus generating more realistic images



uniform



autoregressive

# Generating Diverse High-Fidelity Images with VQ-VAE-2

Ali Razavi et al.

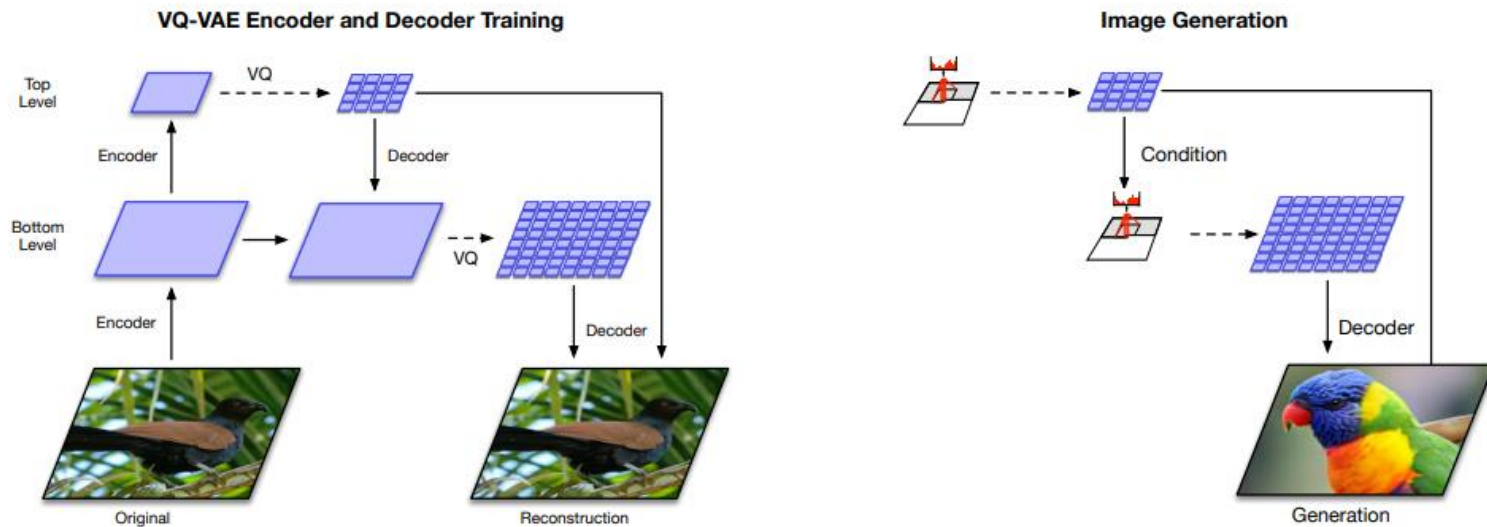
---

Park Jungsoo

Data Mining & Information Systems Lab.  
Department of Computer Science and Engineering,  
College of Informatics, Korea University

---

## Overview



- Top CB for global structure, Bottom CB for the details

## Algorithm

---

**Algorithm 1** VQ-VAE training (stage 1)
 

---

**Require:** Functions  $E_{top}$ ,  $E_{bottom}$ ,  $D$ ,  $\mathbf{x}$   
(batch of training images)

- 1:  $\mathbf{h}_{top} \leftarrow E_{top}(\mathbf{x})$   
      $\triangleright$  quantize with top codebook eq 1
  - 2:  $\mathbf{e}_{top} \leftarrow \text{Quantize}(\mathbf{h}_{top})$
  - 3:  $\mathbf{h}_{bottom} \leftarrow E_{bottom}(\mathbf{x}, \mathbf{e}_{top})$   
      $\triangleright$  quantize with bottom codebook eq 1
  - 4:  $\mathbf{e}_{bottom} \leftarrow \text{Quantize}(\mathbf{h}_{bottom})$
  - 5:  $\hat{\mathbf{x}} \leftarrow D(\mathbf{e}_{top}, \mathbf{e}_{bottom})$   
      $\triangleright$  Loss according to eq 2
  - 6:  $\theta \leftarrow \text{Update}(\mathcal{L}(\mathbf{x}, \hat{\mathbf{x}}))$
- 

---

**Algorithm 2** Prior training (stage 2)
 

---

- 1:  $\mathbf{T}_{top}, \mathbf{T}_{bottom} \leftarrow \emptyset$   $\triangleright$  training set
  - 2: **for**  $\mathbf{x} \in \text{training set}$  **do**
  - 3:    $\mathbf{e}_{top} \leftarrow \text{Quantize}(E_{top}(\mathbf{x}))$
  - 4:    $\mathbf{e}_{bottom} \leftarrow \text{Quantize}(E_{bottom}(\mathbf{x}, \mathbf{e}_{top}))$
  - 5:    $\mathbf{T}_{top} \leftarrow \mathbf{T}_{top} \cup \mathbf{e}_{top}$
  - 6:    $\mathbf{T}_{bottom} \leftarrow \mathbf{T}_{bottom} \cup \mathbf{e}_{bottom}$
  - 7: **end for**
  - 8:  $p_{top} = \text{TrainPixelCNN}(\mathbf{T}_{top})$
  - 9:  $p_{bottom} = \text{TrainCondPixelCNN}(\mathbf{T}_{bottom}, \mathbf{T}_{top})$
  - $\triangleright$  Sampling procedure
  - 10: **while** true **do**
  - 11:    $\mathbf{e}_{top} \sim p_{top}$
  - 12:    $\mathbf{e}_{bottom} \sim p_{bottom}(\mathbf{e}_{top})$
  - 13:    $\mathbf{x} \leftarrow D(\mathbf{e}_{top}, \mathbf{e}_{bottom})$
  - 14: **end while**
-

## Effect of Hierarchical Latent Representation



Figure 3: Reconstructions from a hierarchical VQ-VAE with three latent maps (top, middle, bottom). The rightmost image is the original. Each latent map adds extra detail to the reconstruction. These latent maps are approximately 3072x, 768x, 192x times smaller than the original image (respectively).