Data $x \in \mathbb{R}^D$

Encoder: $z = f_\theta(x)$ : $z \in \mathbb{R}^d$, $d << D$     Bottleneck

Decoder: $\hat{x} = g_\phi(z)$

Training: $\min\limits_{\theta,\phi} E_x[\mathcal{L}_{rec}(x,\hat{x})]$

$z_e \in \mathbb{R}^d$ needs accuracy $\longrightarrow$ +Bits $\propto$ +Computational cost.

$\downarrow$

Solution: Codebook $E = \{e_1, \ldots, e_K\}$ : $e_K \in \mathbb{R}^d$

Instead of giving $z_e$ to $g_\phi(\cdot)$, we quantize it.

$K^*(z) = \arg\min\limits_{K} \|z - e_K\|^2$   (Nearest neighbor) , $z_q = e_K$

$\downarrow$      *Improvement!*

Stores an index, not a specific value. ($\sim \log_2 K$ bits per simbol)

New Problem: The size $K$ of the codebook

• Small $K$: A few bits, but high error $\|z_e - z_q\| \longrightarrow$ worse reconstruction.

• Big $K$: Better accuracy $\|z_e - z_q\| \to 0$, but gigantic codebook $\longrightarrow$ + Computational cost, + Bits.

New Solution! Make a codebook family $\{E^{(m)}\}_{m=1}^{M}$

$E^{(m)} : \{e_1^{(m)}, \ldots, e_K^{(m)}\} \in \mathbb{R}^d$ ; $m = 1, \ldots, M \longrightarrow$ Quantization Levels.

Now, the quantized latent:

$z_q = \sum\limits_{m=1}^{M} q^{(m)}$ ; $q^{(m)} \in E^{(m)}$

The "compressed" $x$ is not a single index, now it is a tuple:

$(k_1, k_2, \ldots, k_M)$ : $k_m \in \{1, \ldots, k_m\}$

And its computational cost is:    $\longrightarrow$ Not necessary every codebook has the same lenght.

$R \tilde{\approx} \sum\limits_{m=1}^{M} \log_2 k_m$

$(k_1^*, \ldots, k_M^*) = \arg\min\limits_{k_1, \ldots, k_M} \left\| z - \sum\limits_{m=1}^{M} e_{k_m}^{(m)} \right\|^2$    Nearest Neighbor for every codebook.

Residual      Quantized