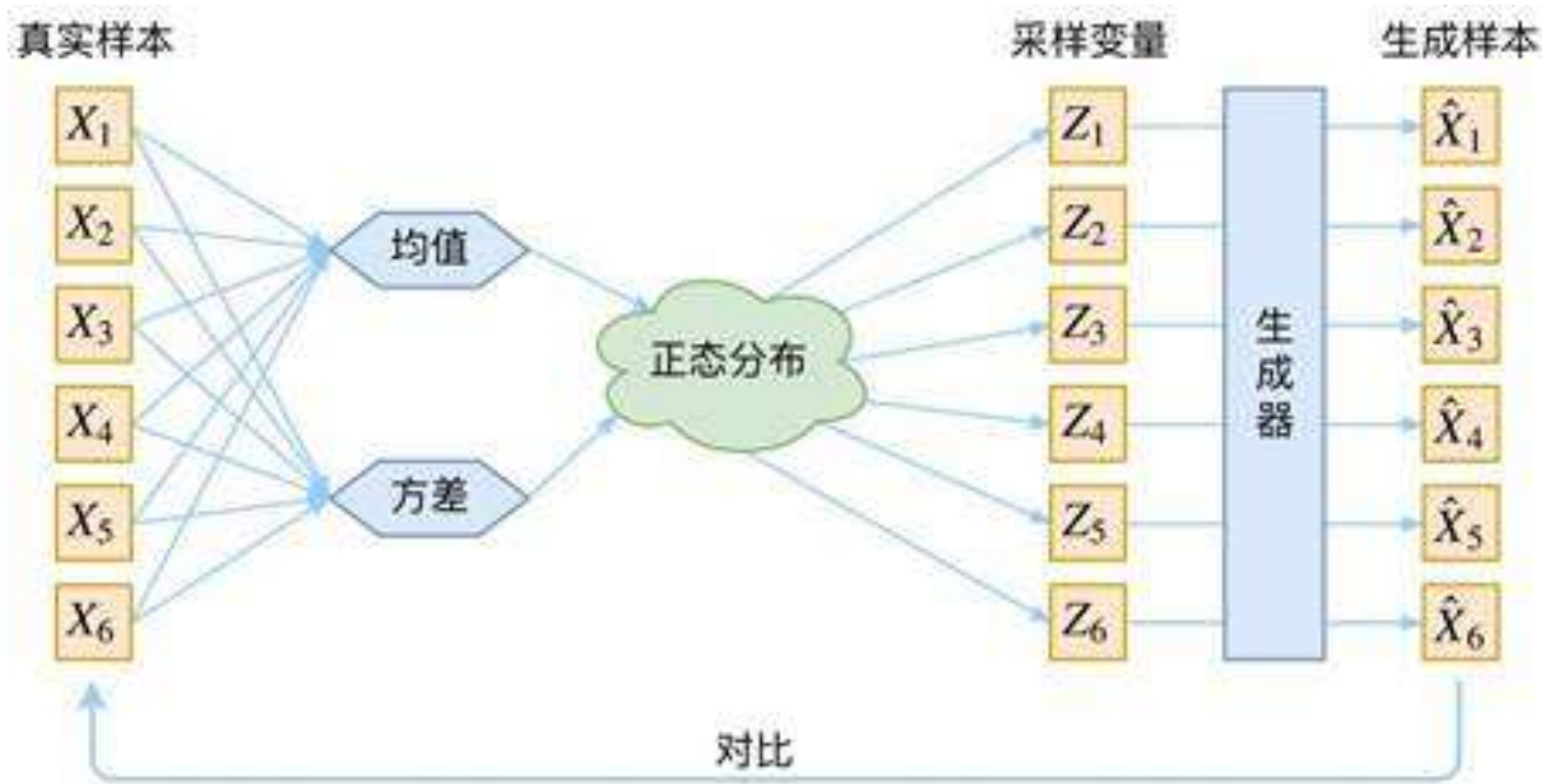


# Neural Discrete Representation Learning (VQ-VAE)

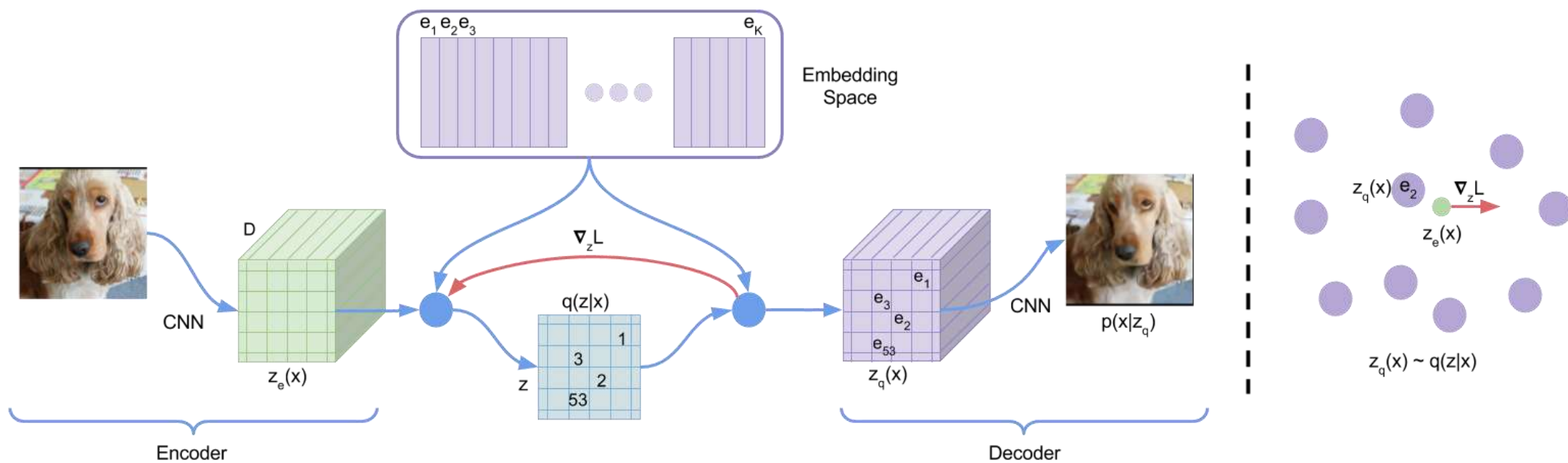
# VAE



# Why do we quantize?

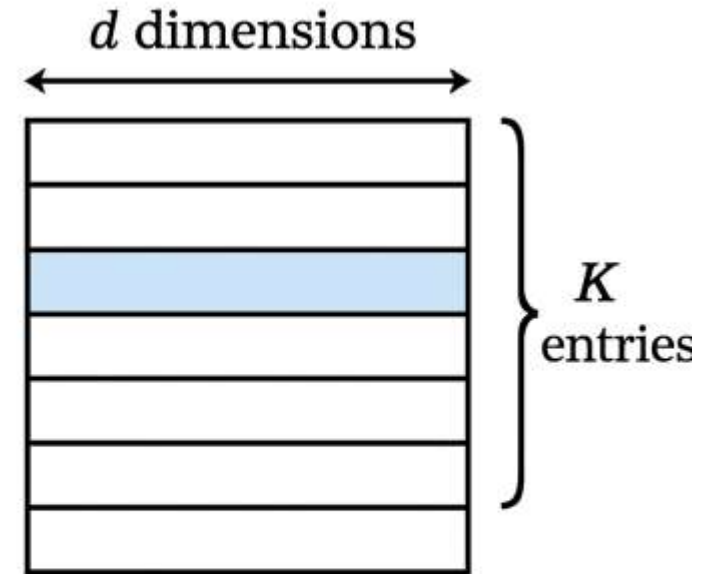
- Discrete representations are potentially a more natural fit for many of the modalities (Audio, Visions and Texts)
- Data compression
- Tokenization

# VQ-VAE



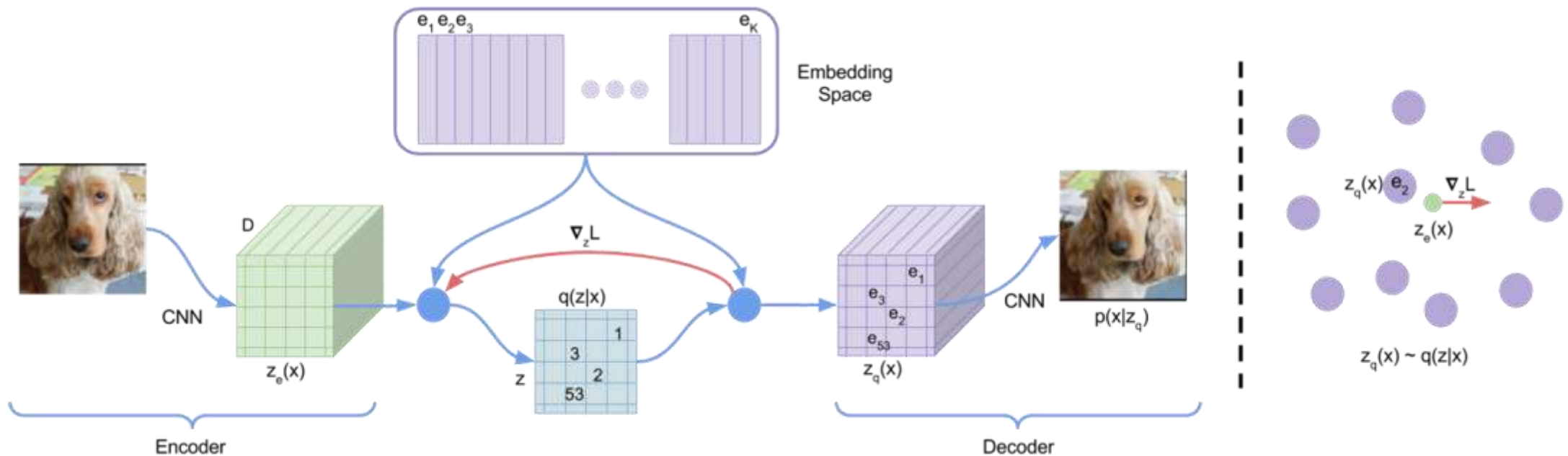
# How do we quantize?

- Map the encoder output  $z_e$  into an entry  $e_i$  of the  $K \times d$  codebook (similar to K-means)
  - Calculate the distance between  $z_e$  and  $e_i$
  - Choose the nearest entry as the input for decoder



**Codebook**

# Straight Through Estimator



# Loss Function

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

Reconstruction Loss      Codebook Loss      Commitment Loss

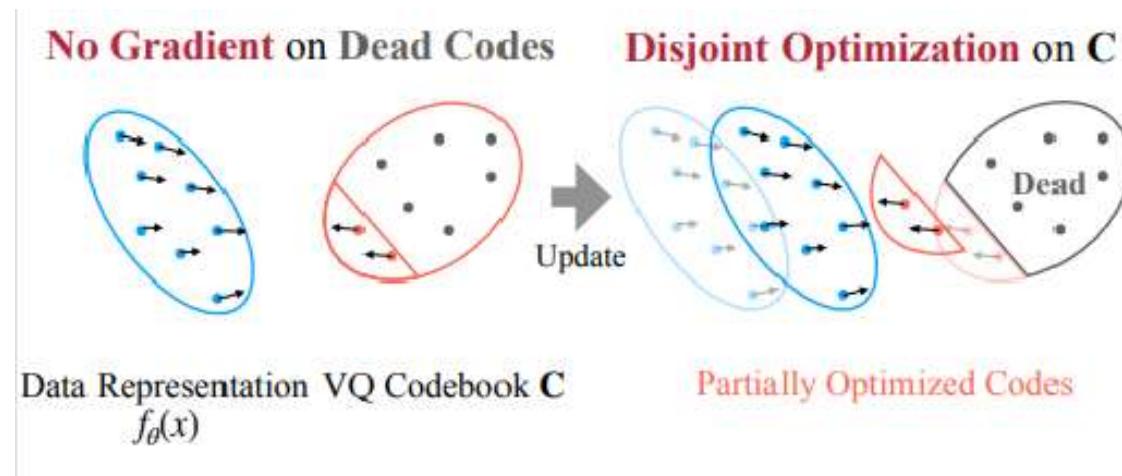
The resulting algorithm to be quite **robust** to  $\beta$ , as the results did not vary for values of  $\beta$  ranging from 0.1 to 2.0.

# Representation Collapse

- For an entry  $e_i$  of the codebook,

$$\frac{\partial L}{\partial e_i} = \begin{cases} \text{not } 0, & \text{if } e_i \text{ is chosen} \\ 0, & \text{if } e_i \text{ is not chosen} \end{cases}$$

- With  $e_i$  and  $z_e$  closer and closer, only a small subset of the codebook entries will be updated.





ADDRESSING REPRESENTATION  
COLLAPSE IN VECTOR  
QUANTIZED MODELS WITH ONE  
LINEAR LAYER (simVQ)

# Representation Collapse

- Representation Collapse:
  - the contradiction between **codebook expansion** and **low codebook utilization** in VQ models where increasing the codebook size fails to improve the performance.
  - the **disjoint optimization process** that updates only a subset of codebook vectors

# SimVQ

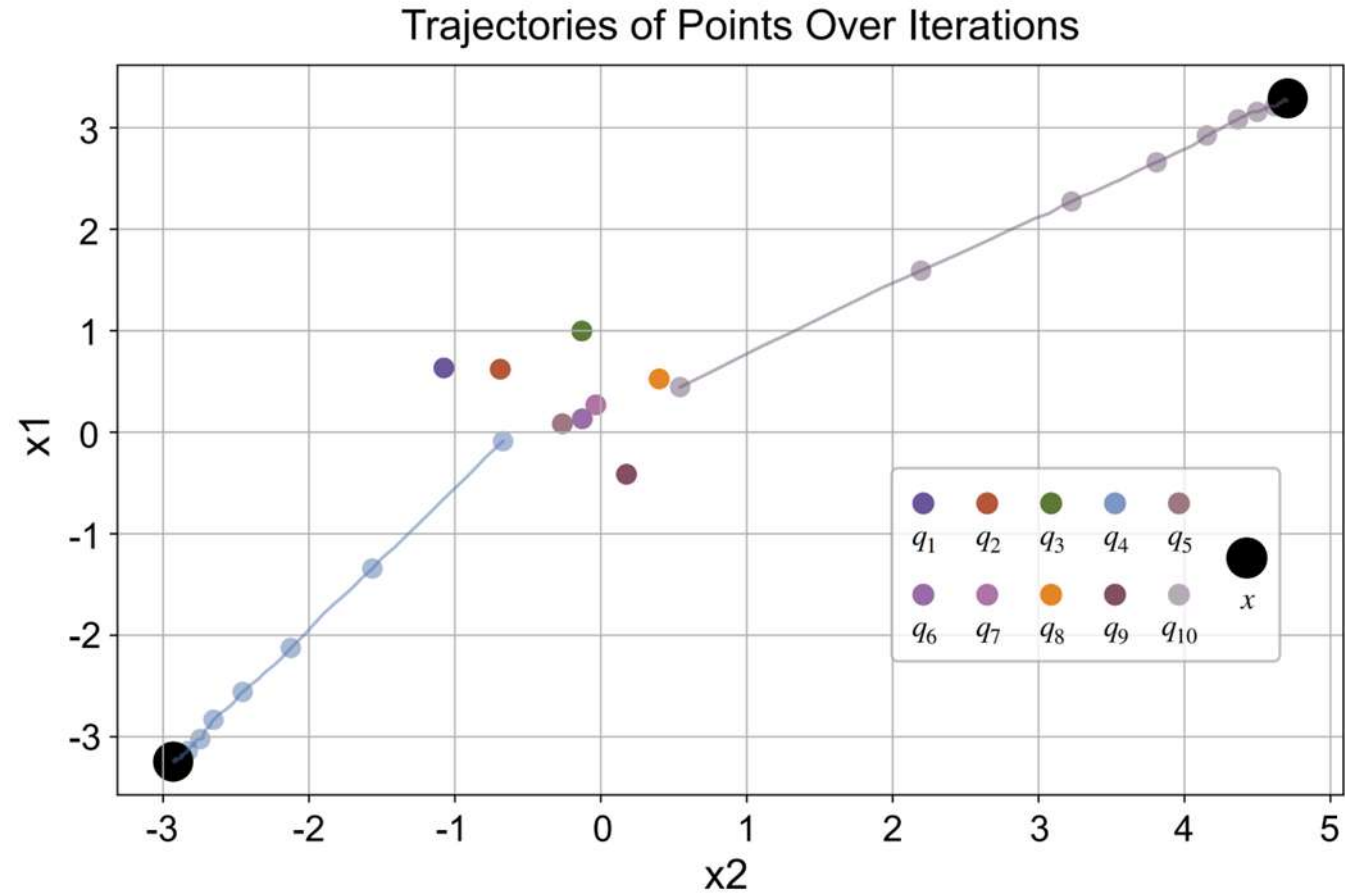
- Reparameterized codebook  $C$  with  $W\pi$ , with  $W = (w_1, w_2, \dots, w_d)$  and  $\pi = (\pi_1, \pi_2, \dots, \pi_K)$ 
  - the optimization of the reparameterized codebook can be divided into three scenarios:
    1. **Updating  $\pi$  with  $W$  frozen:** The vanilla VQ is a special case of this scenario with  $W = I$ .
    2. **Updating  $W$  with  $\pi$  frozen:** The entire codebook  $W\pi$  adjusts to the latent distribution of  $z_e$ . The basis matrix  $W$  rotates and stretches the codebook space.
    3. **Updating both  $W$  and  $\pi$ :** The selected subset of codes moves towards  $z_e$  while the space spanned by  $W$  undergoes simultaneous rotation and stretching.

# Scenario 1: Updating $\boldsymbol{\pi}$ with $\boldsymbol{W}$ frozen

- For one of the entries  $e_i = \boldsymbol{W}\boldsymbol{\pi}_i$ ,

$$\frac{\partial L}{\partial \pi_i} = \frac{\partial L}{\partial e_i} \frac{\partial e_i}{\partial \pi_i} = \frac{\partial L}{\partial e_i} \boldsymbol{W}^T$$

# Scenario 1: Updating $\pi$ with $W$ frozen



## Scenario 2: Updating $W$ with $\pi$ frozen

- For one of the entries  $e_i = W\pi_i$ ,

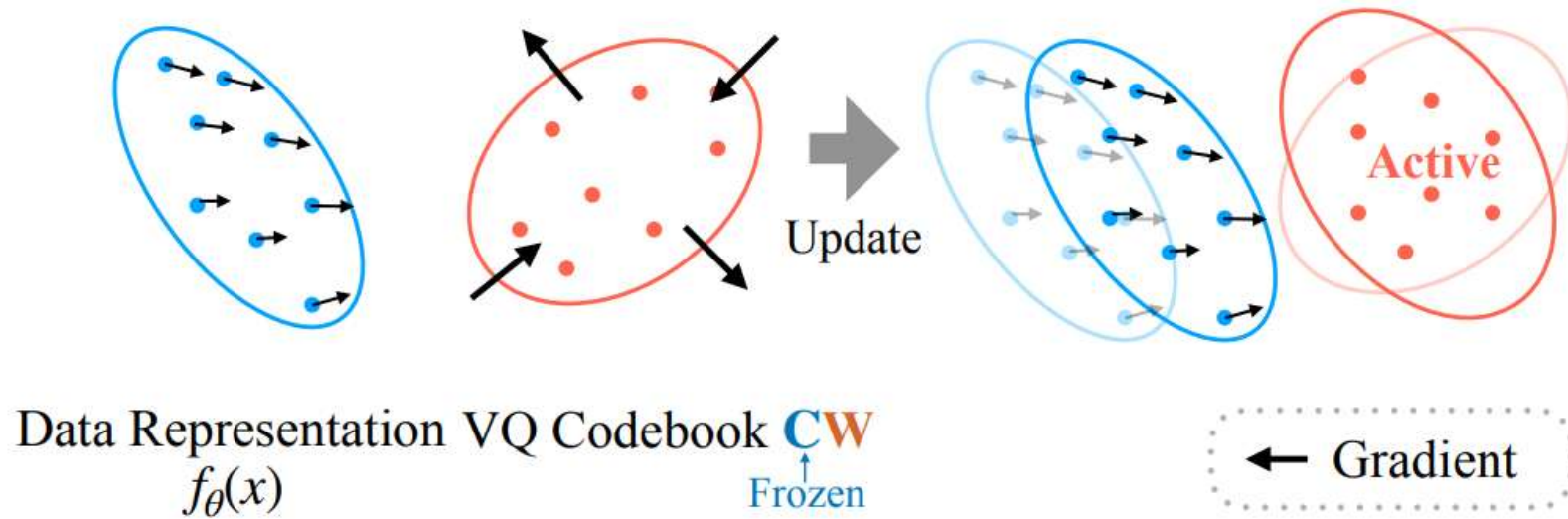
$$\frac{\partial L}{\partial W} = \sum_j \frac{\partial L}{\partial e_j} \frac{\partial e_j}{\partial W} = \sum_j \frac{\partial L}{\partial e_j} \pi_j^T$$

$$e_i^{(t+1)} = W^{(t+1)}\pi_i = \left( W^{(t)} - \eta \frac{\partial L}{\partial W^{(t)}} \right) \pi_i = e_i^{(t)} - \eta \frac{\partial L}{\partial W^{(t)}} \pi_i$$

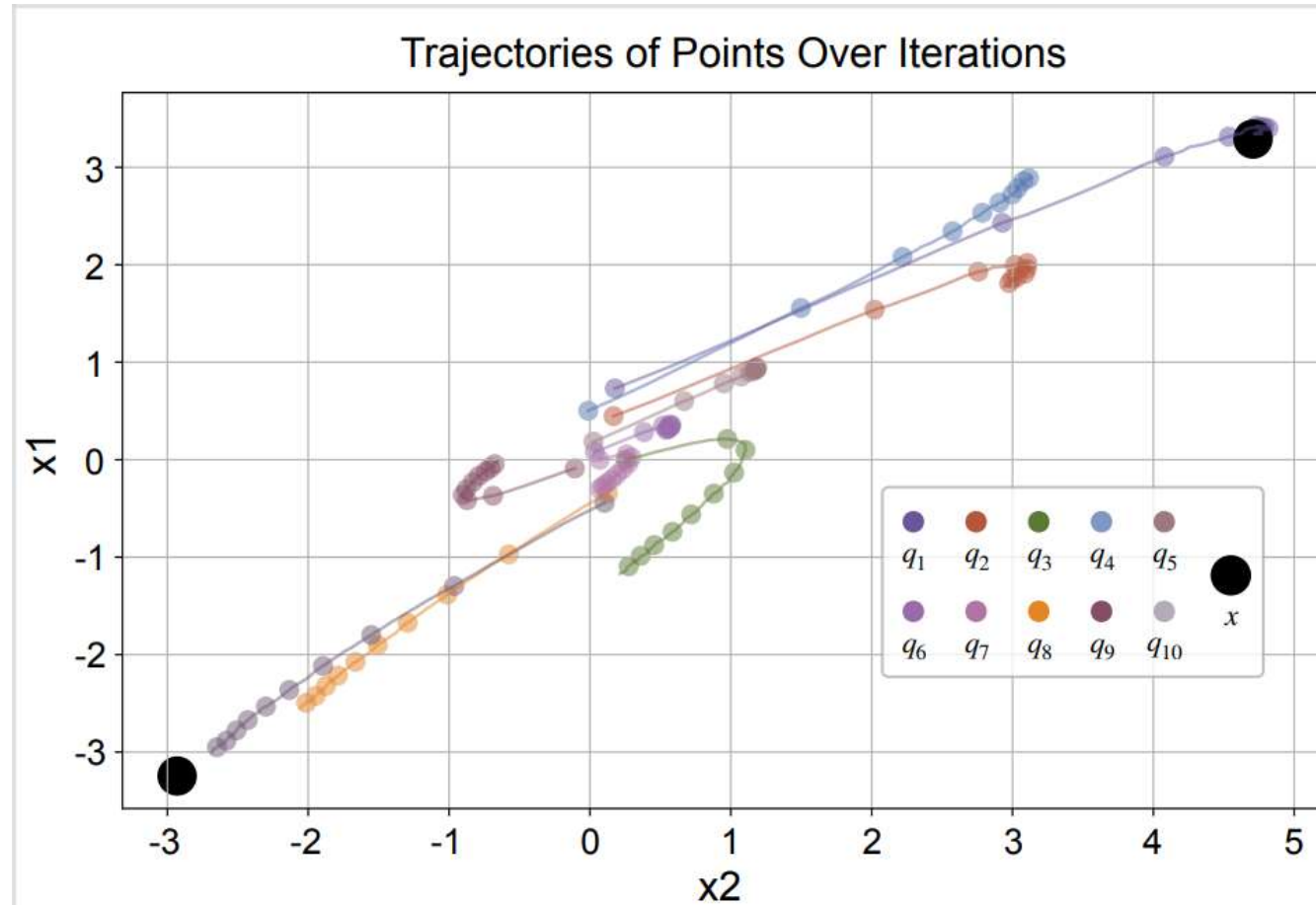
# Scenario 2: Updating $W$ with $\pi$ frozen

## SimVQ

Update **CW** w/ **Latent Basis W**    **All Codes Active** on **CW**

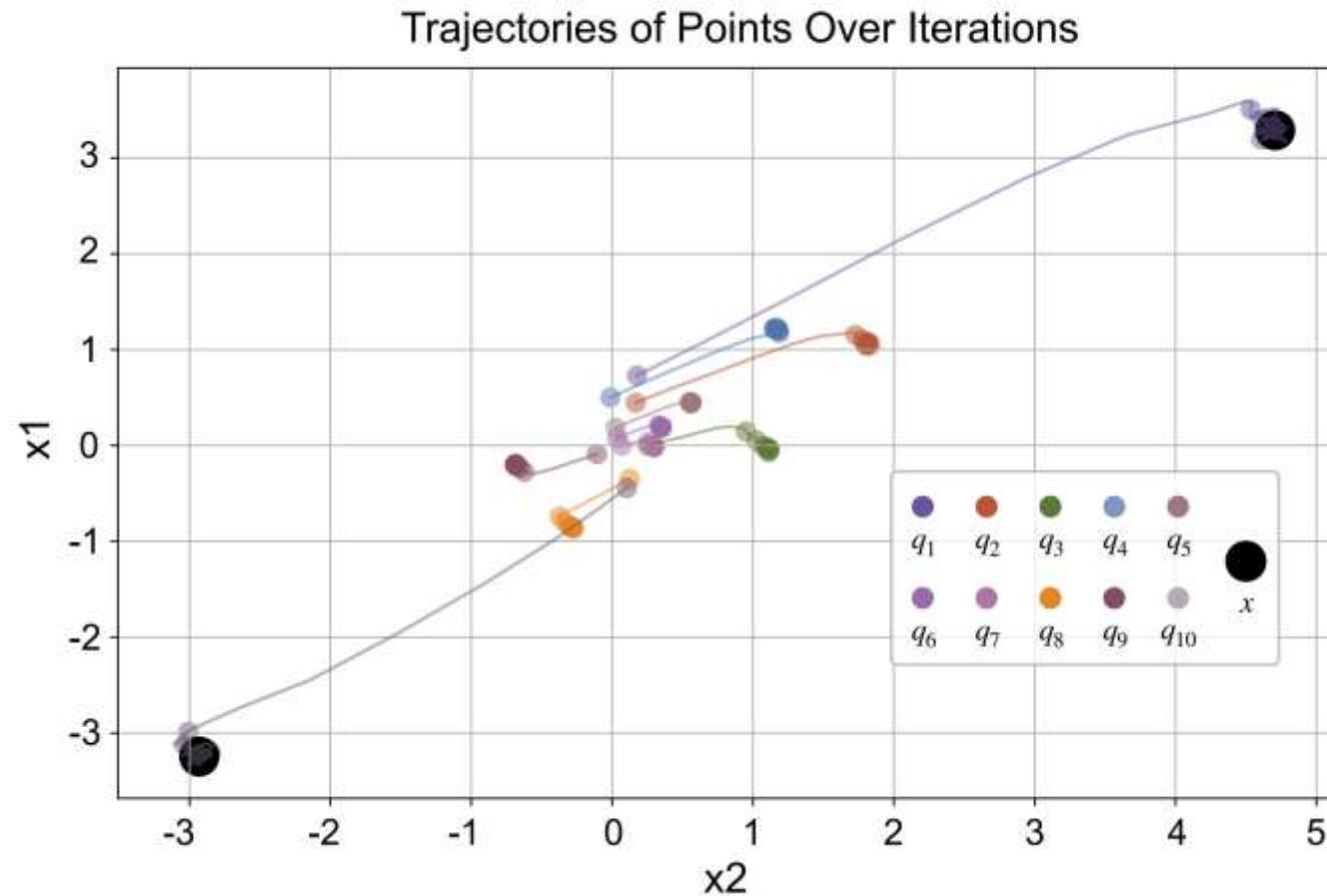


## Scenario 2: Updating $W$ with $\pi$ frozen





# Scenario 3: Updating both $W$ and $\pi$



# Loss Curve



# Experiments

Table 1: Reconstruction performance on ImageNet-1k with a resolution of  $128 \times 128$ . All models are trained using images downsampled into  $16 \times 16$  tokens.  $\dagger$  Results are reproduced using the codebook size of  $[8, 8, 8, 5, 5, 5]$  to approximately match 65,536.  $+$  Following VQGAN-LC, we extract CLIP features with the codebook frozen.

Method	Latent dim	Codebook size	Util $\uparrow$	rFID $\downarrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$
VQGAN (Esser et al., 2021)	128	65,536	1.4%	3.74	0.17	22.20	70.6
VQGAN-EMA (Razavi et al., 2019)	128	65,536	4.5%	3.23	0.15	22.89	72.3
VQGAN-FC (Yu et al., 2022a)	128	65,536	1.4%	5.33	0.18	21.45	68.8
VQGAN-FC (Yu et al., 2022a)	8	65,536	100.0%	2.63	0.13	23.79	77.5
FSQ $^\dagger$ (Mentzer et al., 2024)	16	64,000	100.0%	2.80	0.13	23.63	75.8
LFQ (Yu et al., 2024)	6	65,536	100.0%	2.88	0.13	23.60	77.2
VQGAN-LC-CLIP $^+$ (Zhu et al., 2024a)	768	65,536	100.0%	2.40	0.13	23.98	77.3
SimVQ (ours)	128	65,536	100.0%	<b>2.24</b>	<b>0.12</b>	<b>24.15</b>	<b>78.4</b>
SimVQ (ours)	128	262,144	100.0%	<b>1.99</b>	<b>0.11</b>	<b>24.68</b>	<b>80.3</b>

# Ablation Study on the codebook sizes

Table 2: Ablation study on the effect of various codebook sizes on ImageNet at a resolution of  $128 \times 128$ . † We directly copy the reported results of VQGAN-LC from the original paper on ImageNet  $256 \times 256$  resolution.

Method	Codebook Size	Util↑	rFID↓	LPIPS↓	PSNR↑	SSIM↑
VQGAN-LC-CLIP†	50,000	99.9%	2.75	0.13	23.8	58.4
VQGAN-LC-CLIP†	100,000	99.9%	<u>2.62</u>	0.12	23.8	58.9
VQGAN-LC-CLIP†	200,000	99.8%	<u>2.66</u>	0.12	23.9	59.2
SimVQ	1,024	100.0%	3.67	0.16	22.34	70.8
SimVQ	8,192	100.0%	2.98	0.14	23.23	74.7
SimVQ	65,536	100.0%	2.24	0.12	24.15	78.4
SimVQ	262,144	100.0%	<b>1.99</b>	<b>0.11</b>	<b>24.68</b>	<b>80.3</b>

# Ablation Study on the codebook optimization

Table 3: Ablation study of codebook optimization.

Initialization	Trainable	Util $\uparrow$	rFID $\downarrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$
Gaussian	Yes	100.0%	2.31	0.12	24.04	77.2
Uniform	No	100.0%	2.31	0.12	24.15	78.4
Gaussian	No	100.0%	<b>2.24</b>	<b>0.12</b>	<b>24.15</b>	<b>78.4</b>