

## Motivation

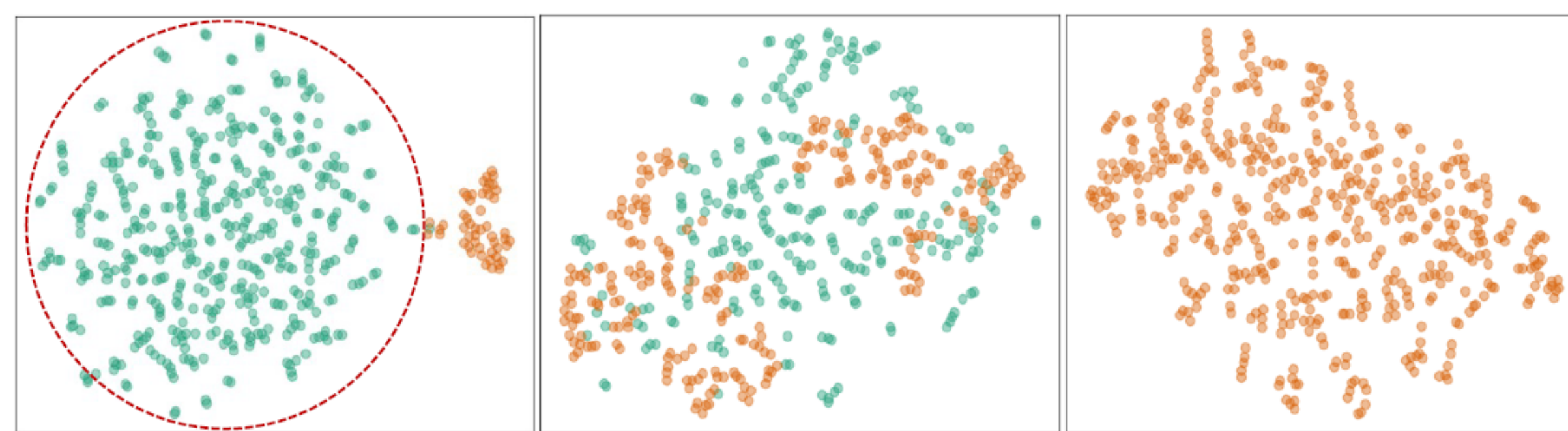
**Goal:** Learn a visual codebook for tasks such as generation with full codevectors utilisation.

**Issue:**

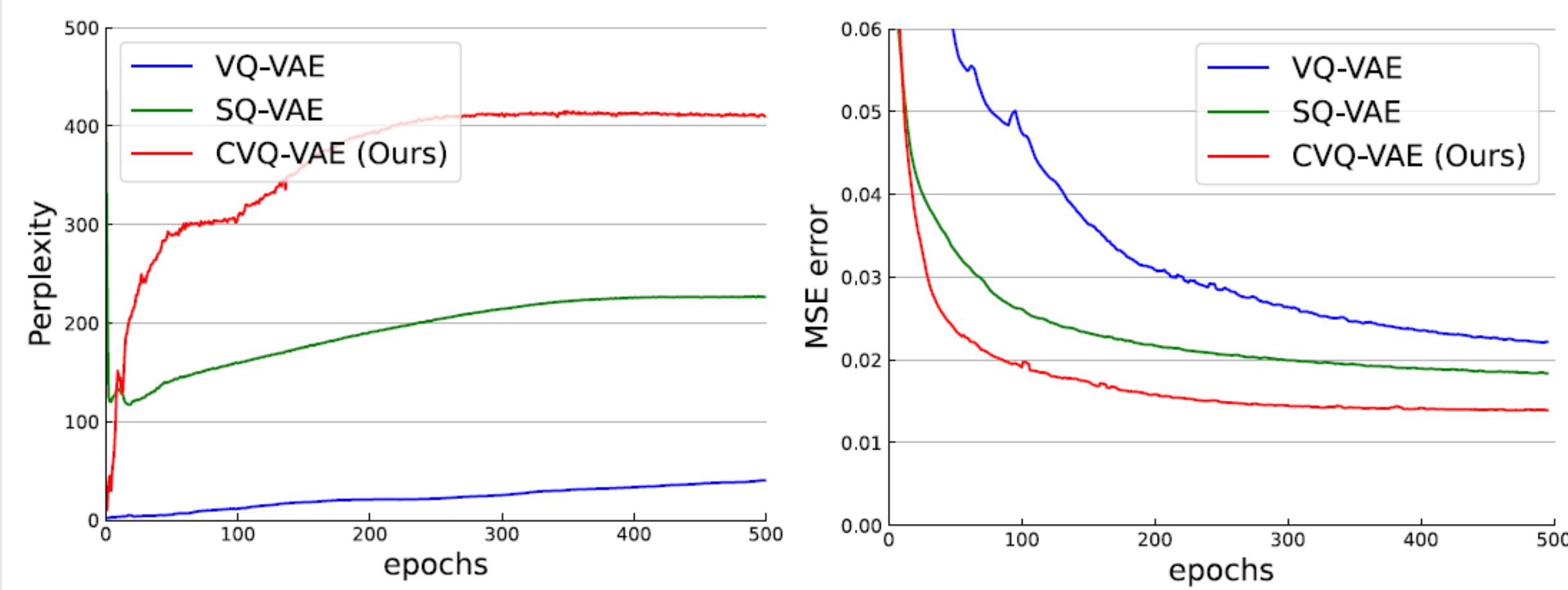
**1. Codebook collapse.** Only a small subset of active codebook entries are optimized

**2. Stop-gradient operator.** Loss can only back propagate to the selected entries.

Green Points: “Dead” Codebook Entries



(a) VQ-VAE [37] Usage: 9.96% (b) SQ-VAE [36] Usage: 49.02% (c) CVQ-VAE Usage: 100%

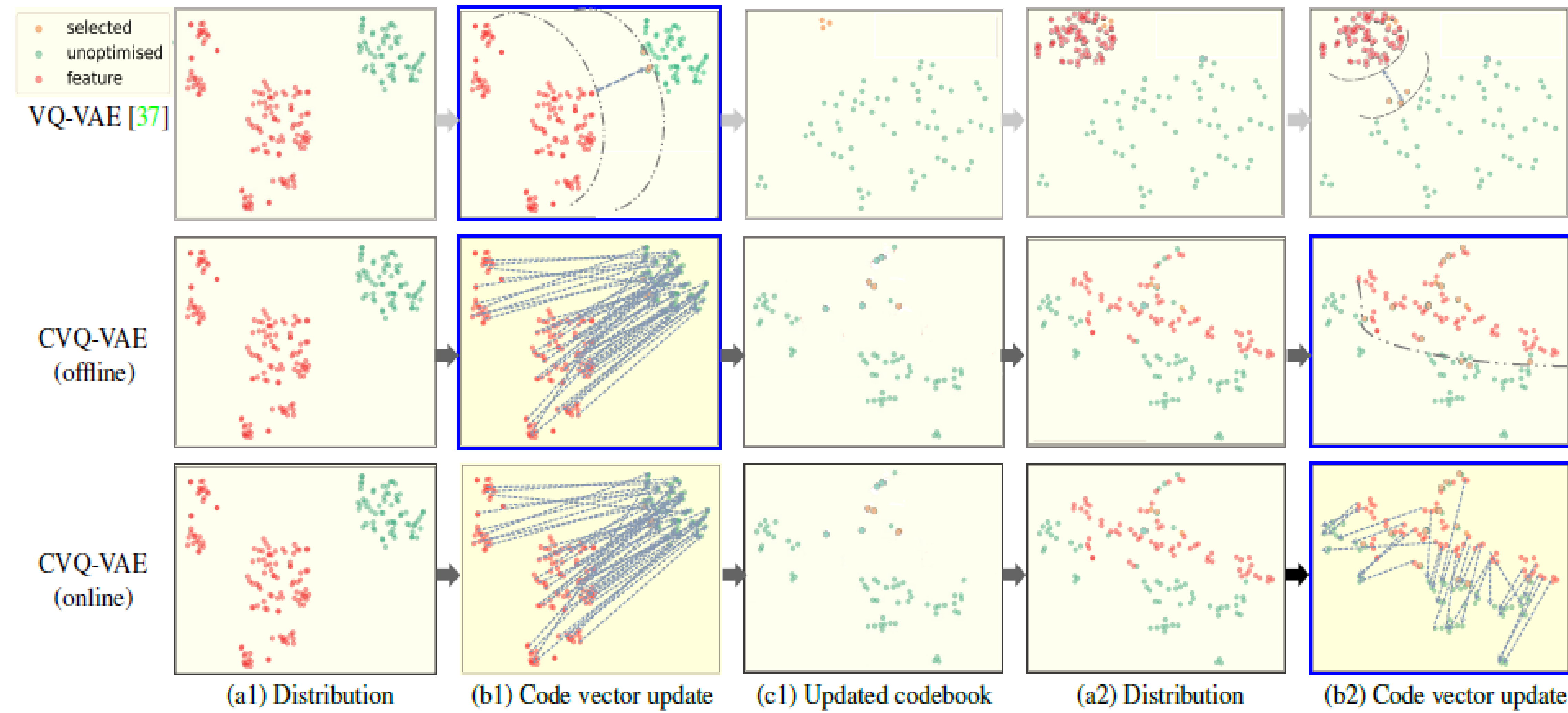


How to ensure the embedded features and codebook entries closely adhere to the same distribution?

## Stage-1: Perceptual Compression



## Contribution: Online Clustered Codebook Optimisation



Codebook Usage:

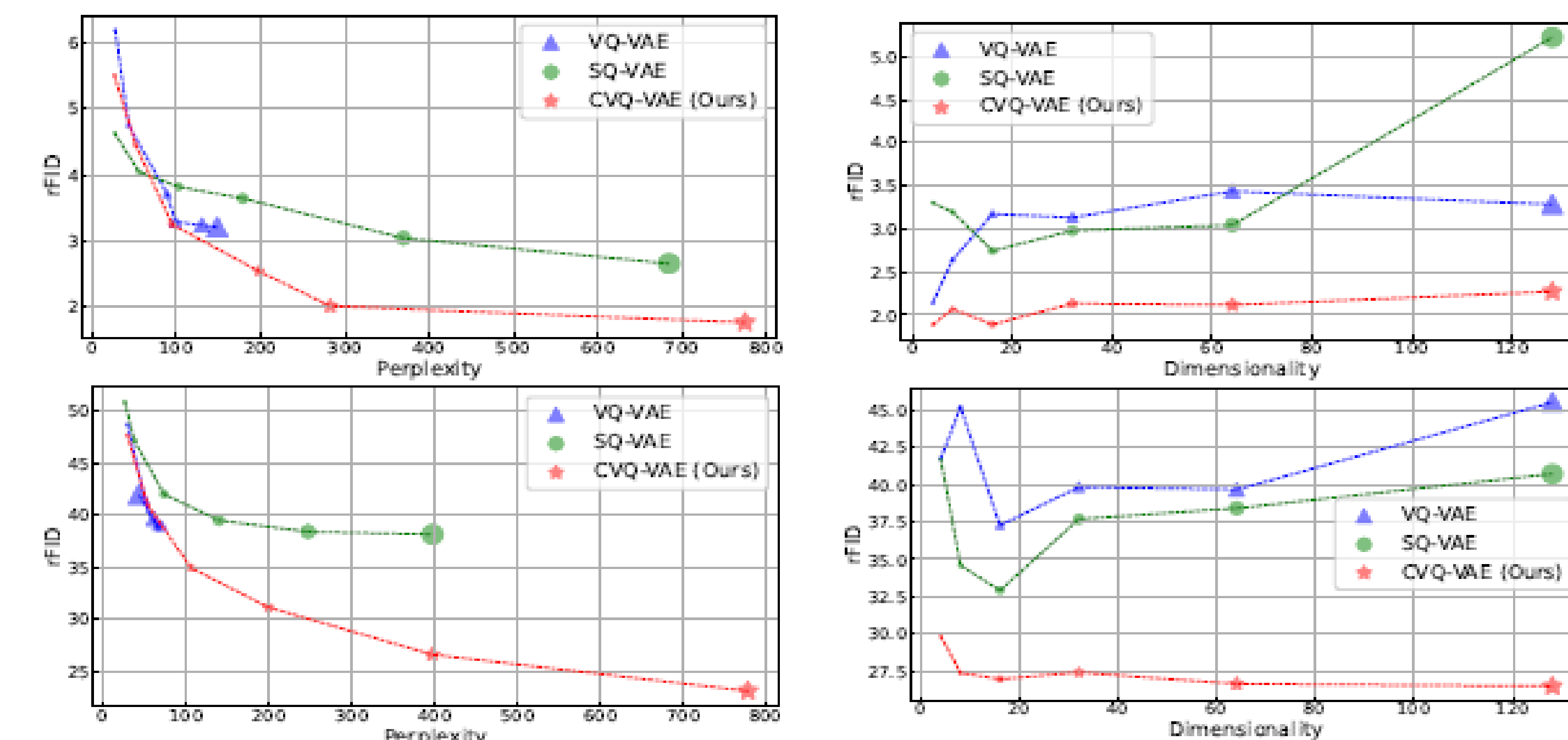
$$N_k^{(t)} = N_k^{(t-1)} \cdot \gamma + \frac{n_k^{(t)}}{Bhw} \cdot (1 - \gamma)$$

Running average updates:

$$\alpha_k^{(t)} = \exp^{-N_k^{(t)} K \frac{10}{1-\gamma} - \epsilon},$$

$$e_k^{(t)} = e_k^{(t-1)} \cdot (1 - \alpha_k^{(t)}) + \hat{z}_k^{(t)} \cdot \alpha_k^{(t)}$$

## Stage-1: Ablation Study on Image Quantization



(a) Codebook size. The blobs' size is proportional to the number of codebook vectors {32, 64, 128, 256, 512, 1024}. The larger size naturally leads to better results in our CVQ-VAE.

(b) Codebook dimensionality. The blob's size refers to the dimensionality of codebook vectors {4, 8, 16, 32, 64, 128}. The higher dimensionality does not ensure a better representation.

(c) Anchor sampling methods. The choice of anchor sampling method has a significant impact on offline (one-time) feature initialization, while the online clustered method is robust for various samplings.

| Methods                   | MNIST (28×28) |               |             | CIFAR10 (32×32) |               |              | FFHQ (256×256) |               |             |
|---------------------------|---------------|---------------|-------------|-----------------|---------------|--------------|----------------|---------------|-------------|
|                           | SSIM ↑        | LPIPS ↓       | rFID ↓      | SSIM ↑          | LPIPS ↓       | rFID ↓       | SSIM ↑         | LPIPS ↓       | rFID ↓      |
| near codevectors [39]     | 0.9790        | 0.0270        | 3.17        | 0.8553          | 0.2553        | 41.08        | 0.7282         | 0.1085        | 4.31        |
| hard encoded features [8] | 0.9814        | 0.0243        | 2.25        | 0.8988          | 0.1978        | 29.16        | 0.7646         | 0.0870        | 3.91        |
| running average (ours)    | <b>0.9823</b> | <b>0.0236</b> | <b>2.23</b> | <b>0.8991</b>   | <b>0.1897</b> | <b>26.62</b> | <b>0.8193</b>  | <b>0.0603</b> | <b>2.94</b> |

(d) Codebook reinitialization methods. In previous works [39, 8], each code entry is associated only with a single feature.

## Stage-2: High-Fidelity Image Generation

### Unconditional Generation on LSUN



### Class-conditional Generation on ImageNet



### Quantitative Results on Image Generation

| Methods                 | FID ↓       |             | Model          | FFHQ  |             | ImageNet |             |
|-------------------------|-------------|-------------|----------------|-------|-------------|----------|-------------|
|                         | Churches    | Bedrooms    |                | Steps | FID ↓       | Steps    | FID ↓       |
| StyleGAN [19]           | 4.21        | 2.35        | RQVAE [22]     | 256   | 10.38       | 1024     | 7.55        |
| DDPM [16]               | 7.89        | 4.90        | MoVQ [44]      | 1024  | 8.52        | 1024     | 7.13        |
| ImageBART [10]          | 7.32        | 5.51        | SQ-VAE [33]    | 200   | 5.17        | 250      | 9.31        |
| Projected-GAN [35]      | <b>1.59</b> | <b>1.52</b> | LDM-4 [31]     | 200   | 4.98        | 250      | 10.56       |
| LDM [32]-8*             | 4.02        | -           | CVQ-VAE (ours) | 200   | <b>4.46</b> | 250      | <b>6.87</b> |
| LDM [32]-4              | -           | 2.95        |                |       |             |          |             |
| LDM [32]-8 (reproduced) | 4.15        | 3.57        |                |       |             |          |             |
| CVQ-VAE-LDM [32]-8      | 3.86        | 3.02        |                |       |             |          |             |

## Other prior related works

- Zheng, C., Vuong, T. L., Cai, J., & Phung, D. [Movq: Modulating quantized vectors for high-fidelity image generation](#). NeurIPS, 2022.. ([MoVQ](#) was reported means a lot to [Kandinsky2.1](#), [Github](#) )
- Vuong, T. L., Le, T., Zhao, H., Zheng, C., Harandi, M., Cai, J., & Phung, D. [Vector Quantized Wasserstein Auto-Encoder](#). ICML 2023.