Motivation
○○○
Token-Based Image Generation
○○○○
Main Observation
○○
Further Verification
○○
Visualizations
○○○

# Image Understanding Makes for A Good Tokenizer for Image Generation

Luting Wang    Yang Zhao[1]    Zijian Zhang[1]    Jiashi Feng[1]    Si Liu[*]    Bingyi Kang[1*]

[1]ByteDance    [*]Corresponding Authors

Friday, December 13

**I₁I ByteDance**

## Table of Contents

Motivation
○●○

Token-Based Image Generation
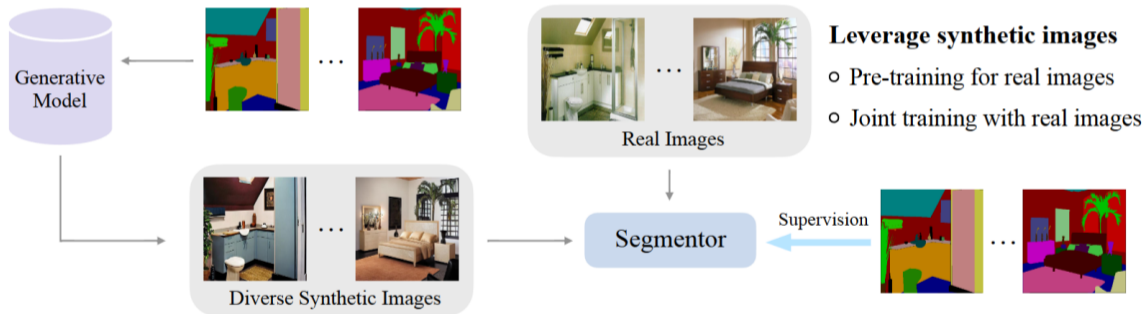○○○○

Main Observation
○○

Further Verification
○○

Visualizations
○○○

# Image Generation Benefits Image Understanding

Studies have shown that IG models can benefit IU tasks in various ways.

**❶ Data augmentation** through synthetic data generation



**Leverage synthetic images**

○ Pre-training for real images

○ Joint training with real images

Real Images

Diverse Synthetic Images

Segmentor

Supervision

Lihe Yang, *et al.* "FreeMask: Synthetic Images with Dense Annotations Make Stronger Segmentation Models." NeurIPS, 2023.

Motivation
○●○

Token-Based Image Generation
○○○○

Main Observation
○○

Further Verification
○○

Visualizations
○○○

# Image Generation Benefits Image Understanding

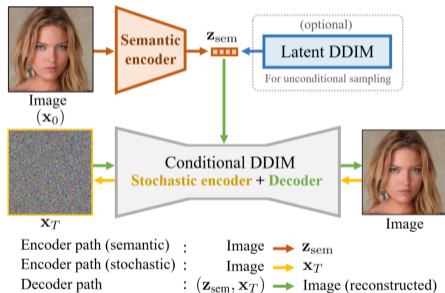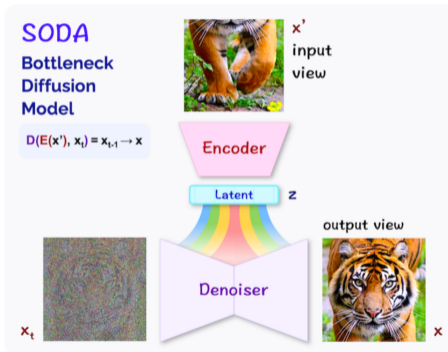Studies have shown that IG models can benefit IU tasks in various ways.
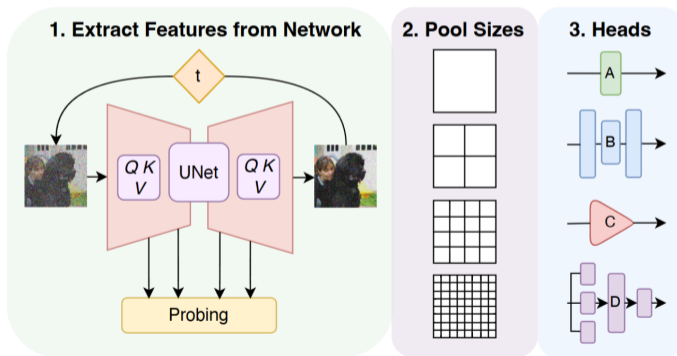
**❷** Improved **representation learning**



Drew A. Hudson, *et al*. "SODA: Bottleneck Diffusion Models for Representation Learning." arXiv:2311.17901, 2023.

Konpat Preechakul, *et al*. "Diffusion Autoencoders: Toward a Meaningful and Decodable Representation." CVPR, 2022.
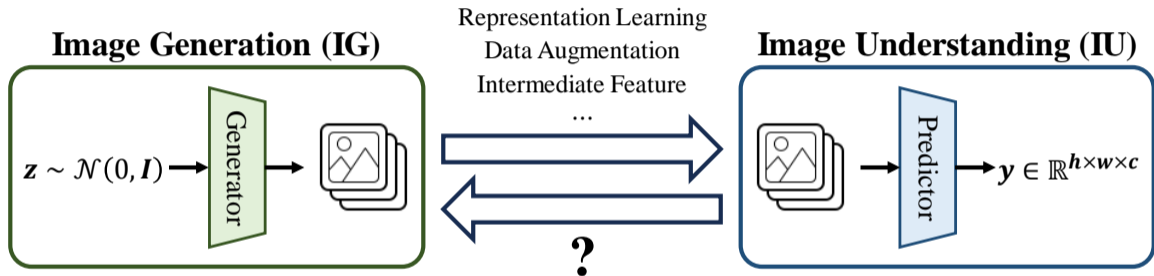
# Image Generation Benefits Image Understanding

Studies have shown that IG models can benefit IU tasks in various ways.

3. Utilizing **intermediate features** from IG models for solving perception tasks



Soumik Mukhopadhyay, *et al.* "Diffusion Models Beat GANs on Image Classification." arXiv:2307.08702.

## The Reciprocal Question?



**Image Generation (IG)**

$z \sim \mathcal{N}(0, I)$ → Generator → 🖼

Representation Learning
Data Augmentation
Intermediate Feature
...

**?**

**Image Understanding (IU)**

🖼 → Predictor → $y \in \mathbb{R}^{h \times w \times c}$

The reciprocal question remains largely uncharted:

**How might IU models aid IG tasks?**

# Table of Contents
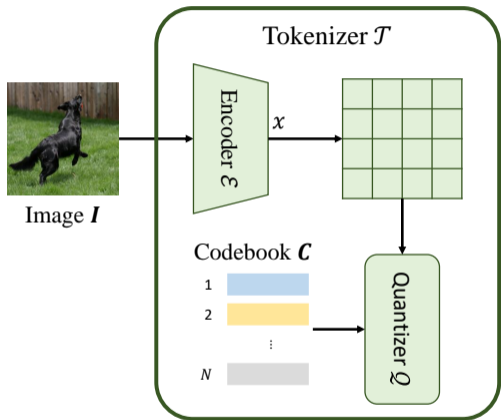
Motivation
○○○

Token-Based Image Generation
○●○○

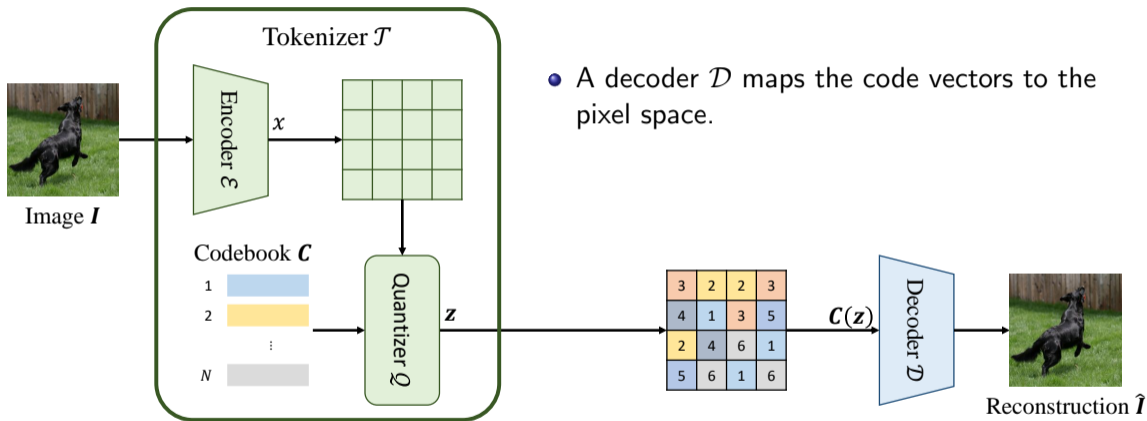Main Observation
○○

Further Verification
○○

Visualizations
○○○

# Two-Stage Image Generation



Tokenizer $\mathcal{T}$

Encoder $\mathcal{E}$

Image $I$

Codebook $C$

Quantizer $\mathcal{Q}$

- The encoder $\mathcal{E}$, quantizer $\mathcal{Q}$, and codebook **C** collectively form an **image tokenizer** $\mathcal{T}$.

- Given an image $I \in \mathbb{R}^{H \times W \times 3}$, the **encoder** $\mathcal{E}$ converts this image into a feature map $x \in \mathbb{R}^{h \times w \times d}$.

- **Codebook C** is a set of $N$ code vectors $\{c_i\}_{i=1}^{N} \in \mathbb{R}^{N \times d}$, where each code vector $c_i \in \mathbb{R}^d$ corresponds to a specific code $i$.

- **Quantizer** $\mathcal{Q}$ then maps $x$ into a sequence of codes $\mathbf{z} = \{z_i\}_{i=1}^{L}$.

Motivation
○○○

Token-Based Image Generation
○●○○

Main Observation
○○

Further Verification
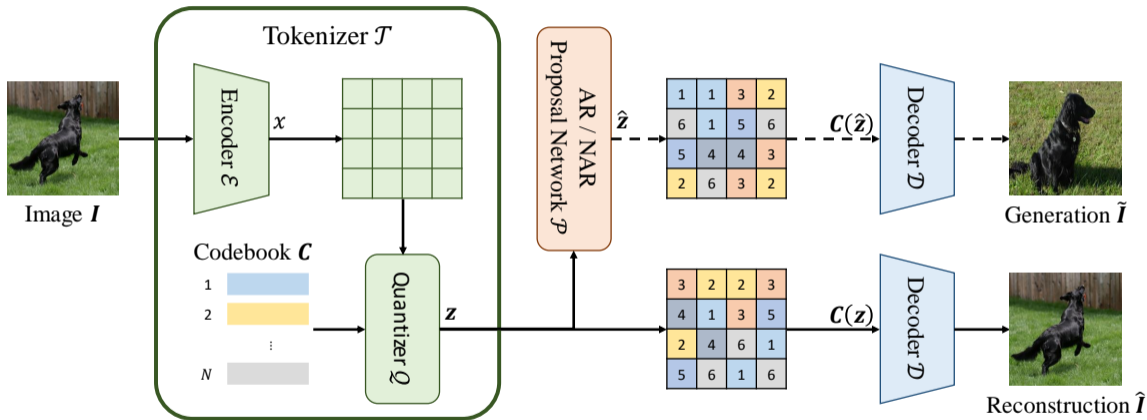○○

Visualizations
○○○

# Two-Stage Image Generation



Image $I$

Tokenizer $\mathcal{T}$

Encoder $\mathcal{E}$

$x$

Codebook $C$

Quantizer $\mathcal{Q}$

$z$

$C(z)$

Decoder $\mathcal{D}$

Reconstruction $\hat{I}$

- A decoder $\mathcal{D}$ maps the code vectors to the pixel space.

Motivation
○○○

Token-Based Image Generation
○●○○

Main Observation
○○

Further Verification
○○

Visualizations
○○○

# Two-Stage Image Generation



- The **proposal network** $\mathcal{P}$ models the distribution over $\mathbf{z}$, denoted as $p(\mathbf{z})$.

Motivation
ooo

Token-Based Image Generation
ooeo

Main Observation
oo

Further Verification
oo

Visualizations
ooo

## Image Tokenizers

Motivation
ooo

Token-Based Image Generation
oooo

Main Observation
oo

Further Verification
oo

Visualizations
ooo

## Image Tokenizers



Vector Quantization

(a) VQGAN

Patrick Esser, Robin Rombach, *et al.* "Taming Transformers for High-Resolution Image Synthesis." CVPR, 2021

Motivation
○○○

Token-Based Image Generation
○○●○

Main Observation
○○

Further Verification
○○

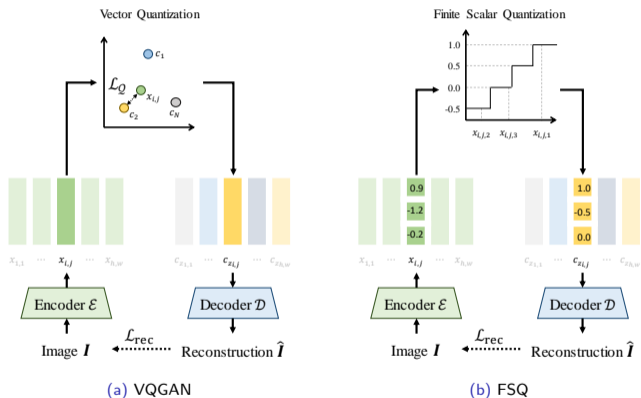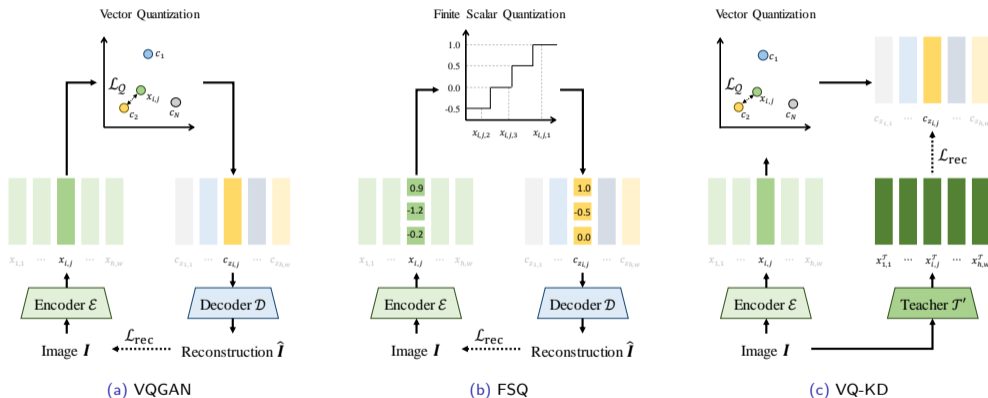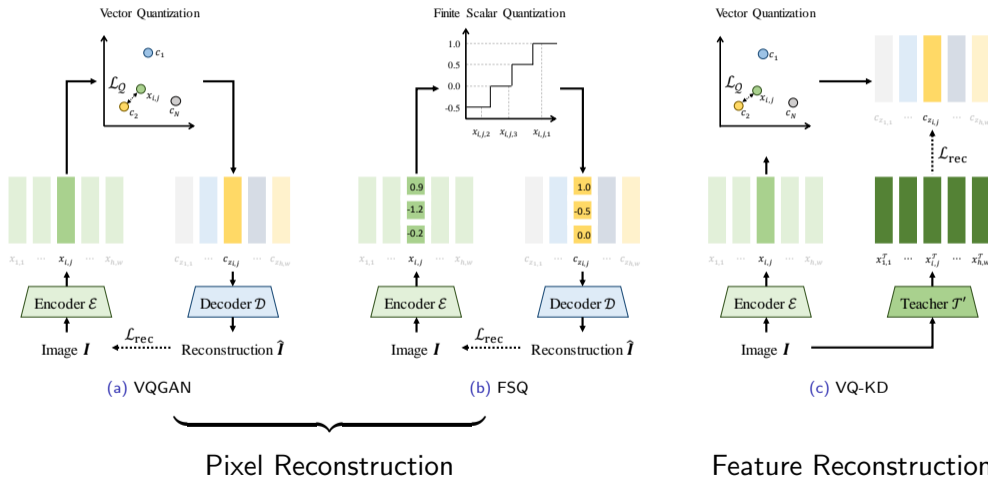Visualizations
○○○

# Image Tokenizers



(a) VQGAN

(b) FSQ

Fabian Mentzer, David Minnen, *et al.* "Finite Scalar Quantization: VQ-VAE Made Simple." ICLR, 2024.

# Image Tokenizers
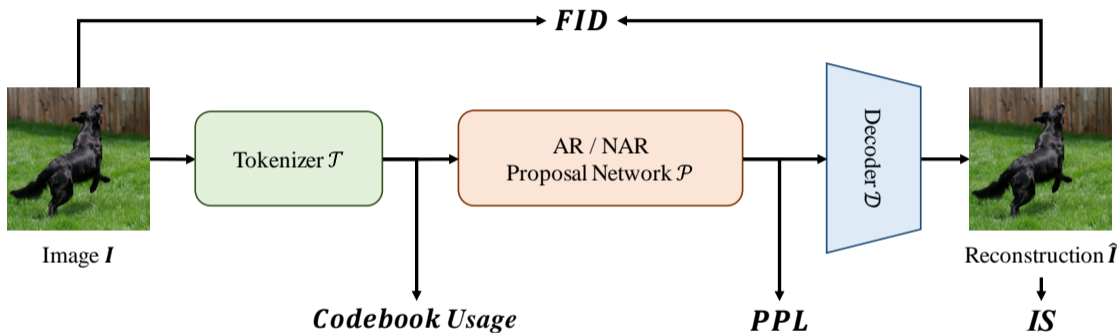


(a) VQGAN

(b) FSQ

(c) VQ-KD

Zhiliang Peng, Li Dong, et al. "BEiT v2: Masked Image Modeling with Vector-Quantized Visual Tokenizers." arXiv preprint arXiv:2208.06366, 2022.

Motivation
ooo

Token-Based Image Generation
oooo

Main Observation
oo

Further Verification
oo

Visualizations
ooo

## Image Tokenizers



(a) VQGAN

(b) FSQ

(c) VQ-KD

Pixel Reconstruction

Feature Reconstruction

Motivation
○○○

Token-Based Image Generation
○○○●

Main Observation
○○

Further Verification
○○

Visualizations
○○○

# Benchmark



We build the above benchmark to evaluate the IG performance of tokenizers.

- For each tokenizer, we train a proposal network and a decoder to form an image generator.
- Various metrics are adopted for a comprehensive evaluation.

Motivation
○○○

Token-Based Image Generation
○○○○

Main Observation
●○

Further Verification
○○

Visualizations
○○○

# Table of Contents

Motivation
ooo

Token-Based Image Generation
oooo

Main Observation
o●

Further Verification
oo

Visualizations
ooo

## Main Observation

1. VQ-KD significantly enhances generation quality over VQGAN.

| Tokenizer | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{NAR}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 4.9 | 5.09 | 116.75 | **24.11** | **20.03** |
| FSQ | 100.0 | 4.96 | 791.56 | 40.17 | 29.78 |
| $VQ\text{-}KD_{CLIP}$ | 100.0 | 4.96 | 53.73 | 11.78 | 9.51 |
| $VQ\text{-}KD_{ViT}$ | 100.0 | 3.69 | 89.30 | **11.40** | **8.45** |
| $VQ\text{-}KD_{DINO}$ | 100.0 | 3.41 | 74.07 | 13.15 | 10.21 |
| $VQ\text{-}KD_{MAE}$ | 100.0 | 4.93 | 280.06 | 26.85 | 16.11 |

Motivation
ooo

Token-Based Image Generation
oooo

Main Observation
o●

Further Verification
oo

Visualizations
ooo

## Main Observation

1. VQ-KD significantly enhances generation quality over VQGAN.
2. The superiority of VQ-KD is irrelevant to the quantization operation and codebook usage.

| Tokenizer | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{NAR}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 4.9 | 5.09 | 116.75 | 24.11 | 20.03 |
| FSQ | **100.0** | **4.96** | 791.56 | **40.17** | 29.78 |
| $VQ\text{-}KD_{CLIP}$ | **100.0** | **4.96** | 53.73 | **11.78** | 9.51 |
| $VQ\text{-}KD_{ViT}$ | 100.0 | 3.69 | 89.30 | 11.40 | 8.45 |
| $VQ\text{-}KD_{DINO}$ | 100.0 | 3.41 | 74.07 | 13.15 | 10.21 |
| $VQ\text{-}KD_{MAE}$ | 100.0 | 4.93 | 280.06 | 26.85 | 16.11 |

Motivation
○○○

Token-Based Image Generation
○○○○

Main Observation
○●

Further Verification
○○

Visualizations
○○○

## Main Observation

1. VQ-KD significantly enhances generation quality over VQGAN.
2. The superiority of VQ-KD is irrelevant to the quantization operation and codebook usage.

| Tokenizer | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{NAR}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 4.9 | 5.09 | **116.75** | 24.11 | 20.03 |
| FSQ | 100.0 | 4.96 | 791.56 | 40.17 | 29.78 |
| $VQ\text{-}KD_{CLIP}$ | 100.0 | 4.96 | **53.73** | 11.78 | 9.51 |
| $VQ\text{-}KD_{ViT}$ | 100.0 | 3.69 | **89.30** | 11.40 | 8.45 |
| $VQ\text{-}KD_{DINO}$ | 100.0 | 3.41 | **74.07** | 13.15 | 10.21 |
| $VQ\text{-}KD_{MAE}$ | 100.0 | 4.93 | 280.06 | 26.85 | 16.11 |

## Main Observation

1. VQ-KD significantly enhances generation quality over VQGAN.
2. The superiority of VQ-KD is irrelevant to the quantization operation and codebook usage.

| Tokenizer | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{NAR}$ ↓ |
|-----------|-------------------|--------|-------|-----------|-------------|
| VQGAN | 4.9 | 5.09 | **116.75** | 24.11 | 20.03 |
| FSQ | 100.0 | 4.96 | **791.56** | 40.17 | 29.78 |
| $VQ\text{-}KD_{CLIP}$ | 100.0 | 4.96 | 53.73 | 11.78 | 9.51 |
| $VQ\text{-}KD_{ViT}$ | 100.0 | 3.69 | 89.30 | 11.40 | 8.45 |
| $VQ\text{-}KD_{DINO}$ | 100.0 | 3.41 | 74.07 | 13.15 | 10.21 |
| $VQ\text{-}KD_{MAE}$ | 100.0 | 4.93 | 280.06 | 26.85 | 16.11 |

Motivation
○○○

Token-Based Image Generation
○○○○

**Main Observation**
○●

Further Verification
○○

Visualizations
○○○

## Main Observation

1. VQ-KD significantly enhances generation quality over VQGAN.
2. The superiority of VQ-KD is irrelevant to the quantization operation and codebook usage.
3. Tokenizers with stronger semantic understanding tend to deliver superior IG performance.

| Tokenizer | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{NAR}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 4.9 | 5.09 | 116.75 | 24.11 | 20.03 |
| FSQ | 100.0 | 4.96 | 791.56 | 40.17 | 29.78 |
| $VQ\text{-}KD_{CLIP}$ | 100.0 | **4.96** | 53.73 | **11.78** | **9.51** |
| $VQ\text{-}KD_{ViT}$ | 100.0 | **3.69** | 89.30 | **11.40** | **8.45** |
| $VQ\text{-}KD_{DINO}$ | 100.0 | **3.41** | 74.07 | **13.15** | **10.21** |
| $VQ\text{-}KD_{MAE}$ | 100.0 | 4.93 | 280.06 | 26.85 | 16.11 |

Motivation
000

Token-Based Image Generation
0000

Main Observation
00

Further Verification
●○

Visualizations
000

# Table of Contents

## Further Verification

1. The superiority of VQ-KD holds across proposal networks.

| Tokenizer | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{NAR}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 4.9 | 5.09 | 116.75 | 24.11 | **20.03** |
| FSQ | 100.0 | 4.96 | 791.56 | 40.17 | **29.78** |
| $VQ\text{-}KD_{CLIP}$ | 100.0 | 4.96 | 53.73 | 11.78 | **9.51** |
| $VQ\text{-}KD_{ViT}$ | 100.0 | 3.69 | 89.30 | 11.40 | **8.45** |
| $VQ\text{-}KD_{DINO}$ | 100.0 | 3.41 | 74.07 | 13.15 | **10.21** |
| $VQ\text{-}KD_{MAE}$ | 100.0 | 4.93 | 280.06 | 26.85 | **16.11** |

## Further Verification

1. The superiority of VQ-KD holds across proposal networks.

| Tokenizer | Codebook Usage (%) | rFID ↓ | PPL ↓ | $\text{FID}_{\text{AR}}$ ↓ | $\text{FID}_{\text{NAR}}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 4.9 | 5.09 | 116.75 | 24.11 | 20.03 |
| FSQ | 100.0 | 4.96 | 791.56 | 40.17 | 29.78 |
| $\text{VQ-KD}_{\text{CLIP}}$ | 100.0 | 4.96 | 53.73 | 11.78 | **9.51** |
| $\text{VQ-KD}_{\text{ViT}}$ | 100.0 | 3.69 | 89.30 | 11.40 | **8.45** |
| $\text{VQ-KD}_{\text{DINO}}$ | 100.0 | 3.41 | 74.07 | 13.15 | 10.21 |
| $\text{VQ-KD}_{\text{MAE}}$ | 100.0 | 4.93 | 280.06 | 26.85 | 16.11 |

## Further Verification

1. The superiority of VQ-KD holds across proposal networks.
2. The superiority of VQ-KD holds across datasets.

| Tokenizer $\mathcal{T}$ | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{T2I}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 2.4 | 16.21 | 47.89 | **38.43** | 24.11 |
| FSQ | 100.0 | 4.62 | 1040.02 | **44.64** | 23.36 |
| VQ-KD$_{CLIP}$ | 82.2 | 5.48 | 72.31 | 29.80 | 11.17 |
| VQ-KD$_{ViT}$ | 100.0 | 3.70 | 117.10 | 23.51 | 15.49 |
| VQ-KD$_{DINO}$ | 100.0 | 2.69 | 129.93 | **17.55** | 11.50 |
| VQ-KD$_{MAE}$ | 100.0 | 3.51 | 317.98 | 44.01 | 15.60 |

Motivation
○○○

Token-Based Image Generation
○○○○

Main Observation
○○

Further Verification
○●

Visualizations
○○○

## Further Verification

1. The superiority of VQ-KD holds across proposal networks.
2. The superiority of VQ-KD holds across datasets.

| Tokenizer $\mathcal{T}$ | Codebook Usage (%) | rFID ↓ | PPL ↓ | $\text{FID}_{AR}$ ↓ | $\text{FID}_{T2I}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | **2.4** | 16.21 | 47.89 | 38.43 | 24.11 |
| FSQ | 100.0 | 4.62 | 1040.02 | 44.64 | 23.36 |
| VQ-KD$_{CLIP}$ | 82.2 | 5.48 | 72.31 | 29.80 | 11.17 |
| VQ-KD$_{ViT}$ | 100.0 | **3.70** | 117.10 | **23.51** | 15.49 |
| VQ-KD$_{DINO}$ | 100.0 | <span style="color:red">**2.69**</span> | 129.93 | <span style="color:red">**17.55**</span> | 11.50 |
| VQ-KD$_{MAE}$ | 100.0 | 3.51 | 317.98 | 44.01 | 15.60 |

# Further Verification

1. The superiority of VQ-KD holds across proposal networks.
2. The superiority of VQ-KD holds across datasets.
3. The superiority of VQ-KD holds across tasks.

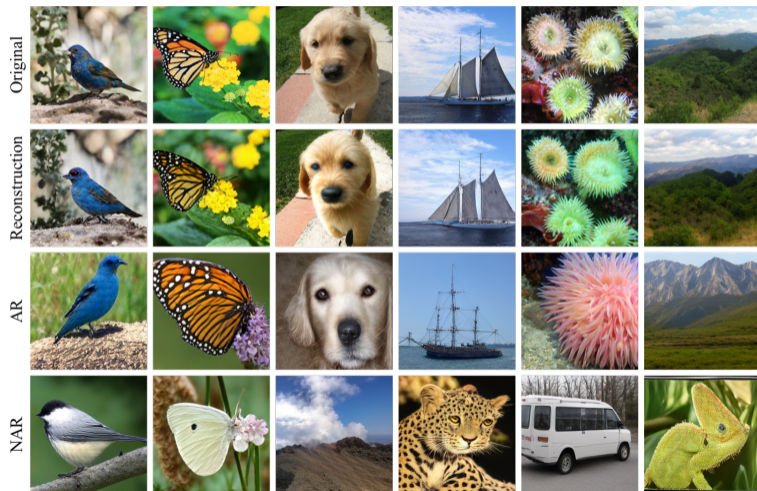| Tokenizer $\mathcal{T}$ | Codebook Usage (%) | rFID ↓ | PPL ↓ | $FID_{AR}$ ↓ | $FID_{T2I}$ ↓ |
|---|---|---|---|---|---|
| VQGAN | 2.4 | 16.21 | 47.89 | 38.43 | **24.11** |
| FSQ | 100.0 | 4.62 | 1040.02 | 44.64 | **23.36** |
| VQ-KD$_{CLIP}$ | 82.2 | 5.48 | 72.31 | 29.80 | **11.17** |
| VQ-KD$_{ViT}$ | 100.0 | 3.70 | 117.10 | 23.51 | **15.49** |
| VQ-KD$_{DINO}$ | 100.0 | 2.69 | 129.93 | 17.55 | **11.50** |
| VQ-KD$_{MAE}$ | 100.0 | 3.51 | 317.98 | 44.01 | **15.60** |

# Table of Contents

Motivation
ooo
Token-Based Image Generation
oooo
Main Observation
oo
Further Verification
oo
Visualizations
o●o

# Results
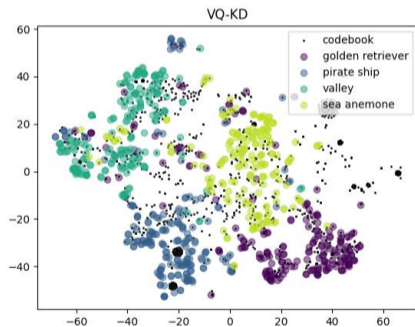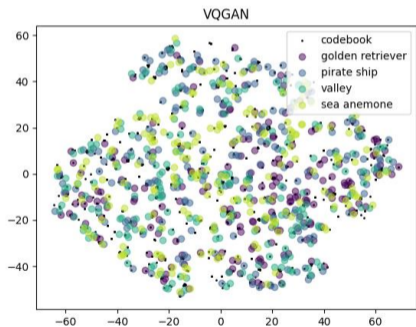


VQ-KD visualization of

- the original images,
- The reconstructed images,
- The AR generation,
- The NAR generation.

Motivation
000

Token-Based Image Generation
0000

Main Observation
00

Further Verification
00

Visualizations
00●

## Codebook

- Compared to VQGAN, the organized feature space of VQ-KD improves the **clarity of code semantics** and helps to **better understand** image content and code interaction.



Codebook visualization of VQGAN and VQ-KD$_{ViT}$.

# **Contact Us**

📄 arxiv.org/abs/2411.04406          @ lutingwang.ai@qq.com

⦿ https://github.com/magic-research/vector_quantization

**Il ByteDance**