# Lab5 - MaskGIT for Image Inpainting

2024 Spring

詹雨婷

# Important Date
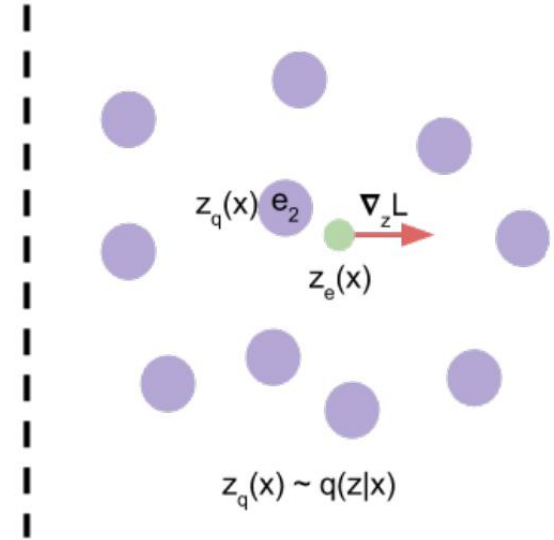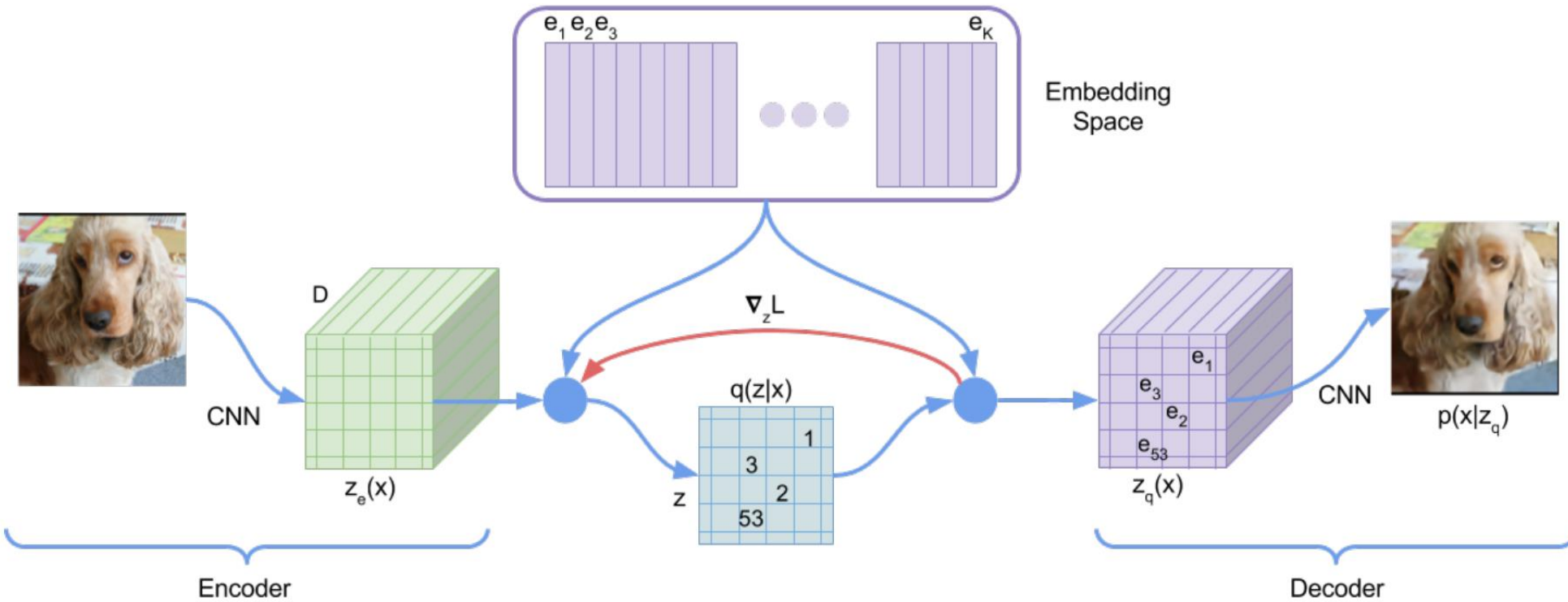
| | LAB1<br>Back-Propagation | LAB2<br>CNN | LAB3<br>CNN | LAB4<br>RNN+VAE | LAB5<br>MaskGIT | LAB6<br>Generative Models |
|---|---|---|---|---|---|---|
| Announce | 3/12 (Tabc) | 3/26 (Tabc) | 4/2 (Tabc) | 4/11 (Rn56) | 4/30 (Tabc) | 5/21 (Tabc) |
| DEMO | 3/26 (Tabc) | 4/11 (Rn56) | 4/11 (Rn56) | 5/7 (Tabc) | 5/21(Rn56) | No demo |

# Submission

- Score: 70% demo score + 40% report

- If the zip file name or the report spec have format error, you will be punished (-5)

- Submission Deadline: 5/21 (Tue) 11:59 a.m.

- Demo date: 5/21 (Tue)

- Turn in: a. Experiment Report (.pdf) b. Source code

- Notice : zip all files in one file and name it like「DL_LAB5_YourStudentID_ name.zip」, ex: [DL_LAB5_312581028_詹雨婷.zip」
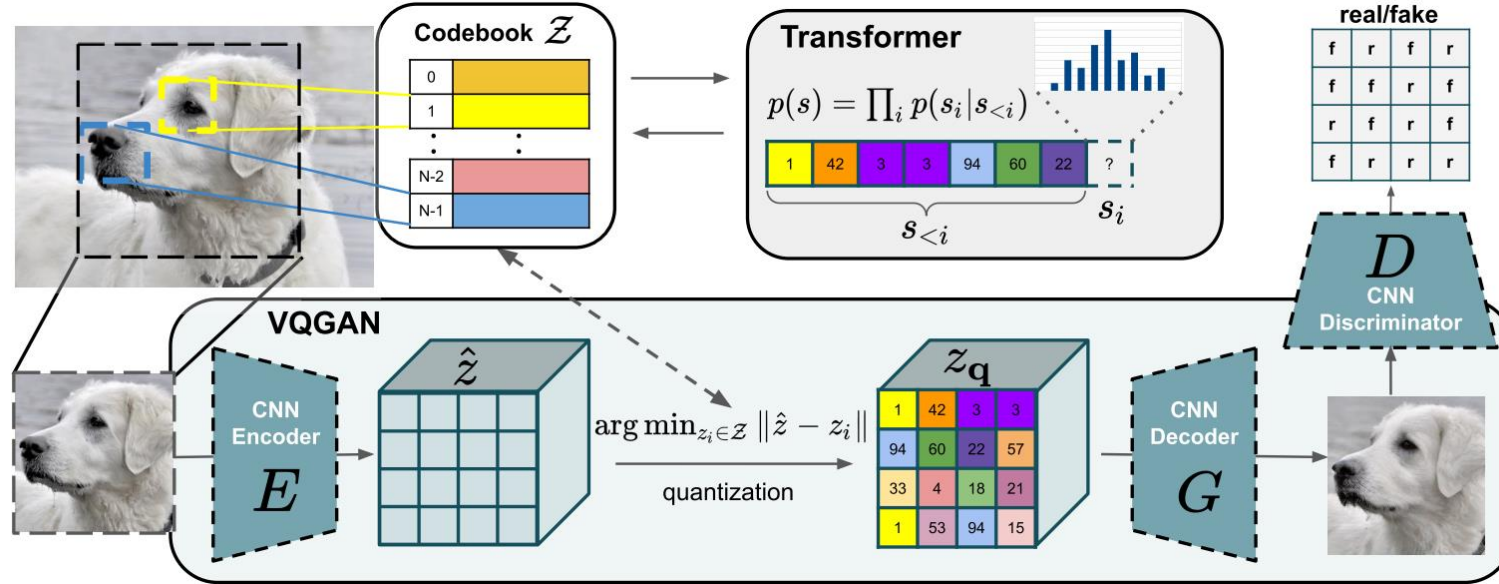
# Introduction

# VQ-VAE (prior work)



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

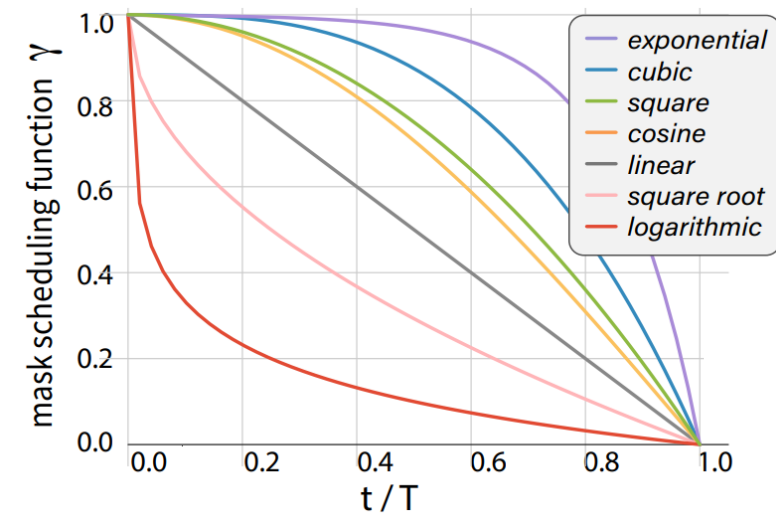- PixelCNN (AR model) prior ancestral sampling z

# VQ-GAN (prior work)



$$\mathcal{L}_{\mathrm{VQ}}(E, G, \mathcal{Z}) = \boxed{\|x - \hat{x}\|^2} + \|\mathrm{sg}[E(x)] - z_{\mathbf{q}}\|_2^2 + \|\mathrm{sg}[z_{\mathbf{q}}] - E(x)\|_2^2.$$

- Perceptual loss replace L2 loss

$$\mathcal{L}_{\mathrm{GAN}}(\{E, G, \mathcal{Z}\}, D) = [\log D(x) + \log(1 - D(\hat{x}))]$$

- Transformer (AR model) prior ancestral sampling z
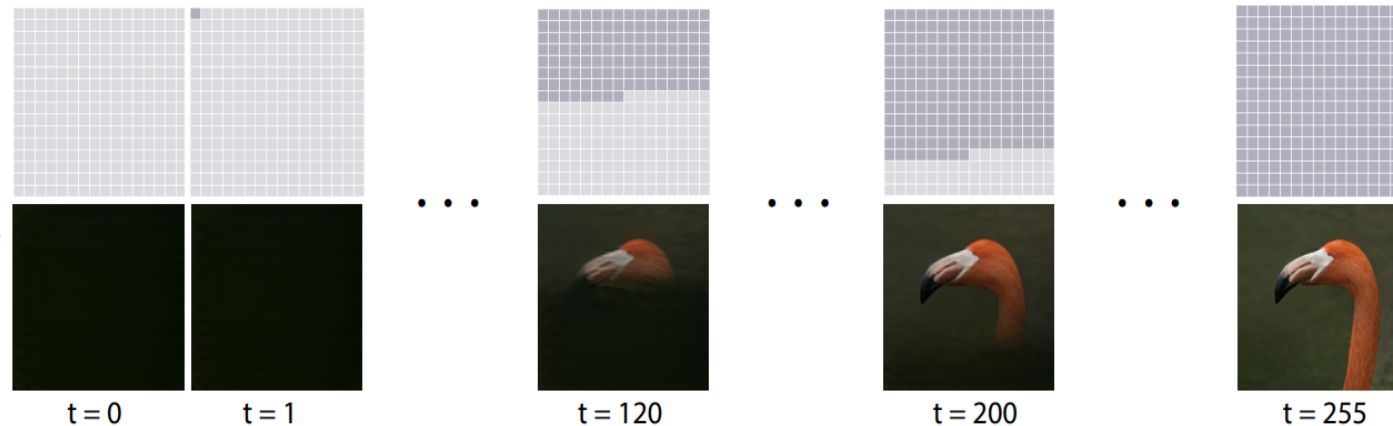
# MaskGIT Pipeline Overview

**STAGE1** Tokenization

**STAGE2** Masked Visual Token Modeling (MVTM)



- Transformer (BERT) prior ancestral sampling z
- MVTM in Training $\gamma(r) \in (0,1]$ $\mathcal{L}_{\text{mask}} = -\mathop{\mathbb{E}}_{\mathbf{Y} \in \mathcal{D}} \left[ \sum_{\forall i \in [1,N], m_i = 1} \log p(y_i | Y_{\overline{\mathbf{M}}}) \right]$
- Iterative Decoding

$$n = \lceil \gamma(\tfrac{t}{T}) N \rceil$$

$$m_i^{(t+1)} = \begin{cases} 1, & \text{if } c_i < \text{sorted}_j(c_j)[n]. \\ 0, & \text{otherwise.} \end{cases}$$

# Iterative Decoding

**VQGAN**

Sequential Decoding with Autoregressive Transformers

| | | | | |
|---|---|---|---|---|
| t = 0 | t = 1 | t = 120 | t = 200 | t = 255 |

**MaskGIT**

Scheduled Parallel Decoding with MaskGIT

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| t = 0 | t = 1 | t = 2 | t = 3 | t = 4 | t = 5 | t = 6 | t = 7 |

# Lab Details

# Lab Objective

- Focus on implementing MaskGIT for the inpainting task

- During testing, images contain gray regions indicating missing information, which we aim to restore using MaskGIT.

- The key practical emphasis of this lab lies in three main areas:
  - Multi-head attention
  - Transformer training
  - Inference inpainting

# Dataset

**a.** **Training dataset**

image: 12000 png files **(./cat_face/train)**

**b.** **Validation dataset**

image: 3000 png files  **(./cat_face/val)**

**c.** **Testing dataset**

masked image: 747 png files **(./cat_face/masked_image)**

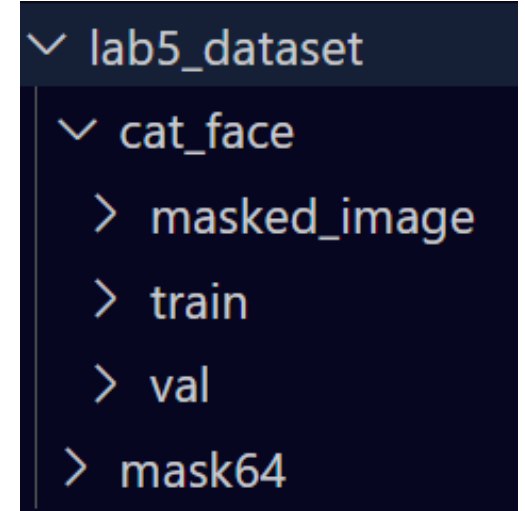mask: 747 png files **(./mask64 )**

**d.** **Download dataset**

i. ON your own machine

?> sftp -P 10046 pp037@140.113.215.196 (passwd: pp037OnClass)

?> get lab5_dataset.zip

ii. ON Provided machine

?> sftp pp037@192.168.201.46 (passwd: pp037OnClass)

?> get lab5_dataset.zip

Reference: https://www.kaggle.com/datasets/spandan2/cats-faces-64x64-for-generative-models

lab5_dataset
cat_face
masked_image
train
val
mask64

# VQGAN Stage1 Pretrained Weight

- **You can't modify any model structure or retrain stage1.**
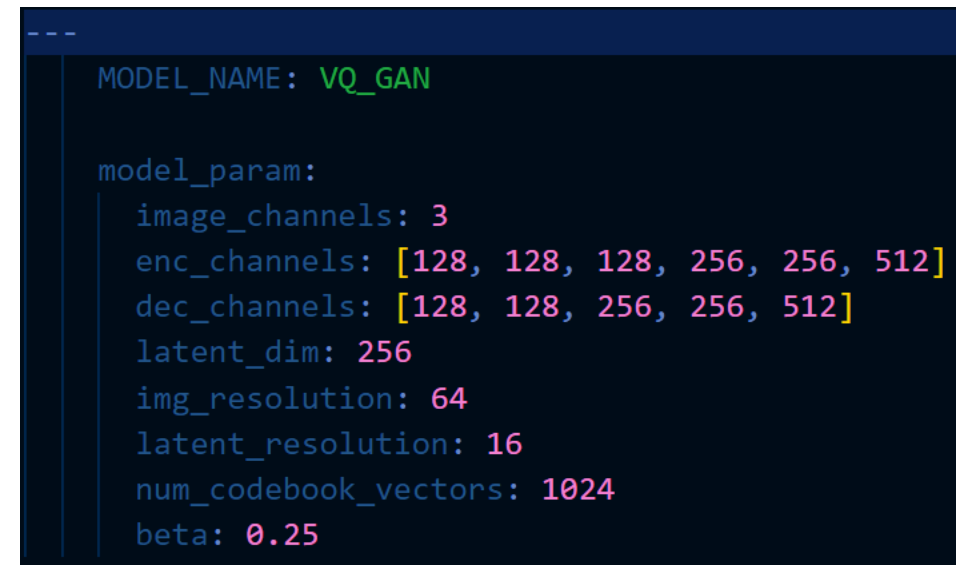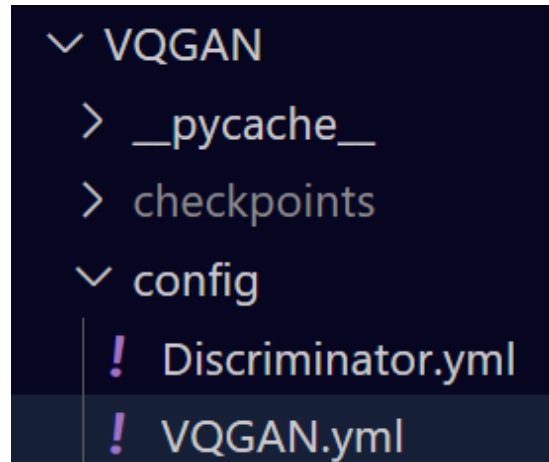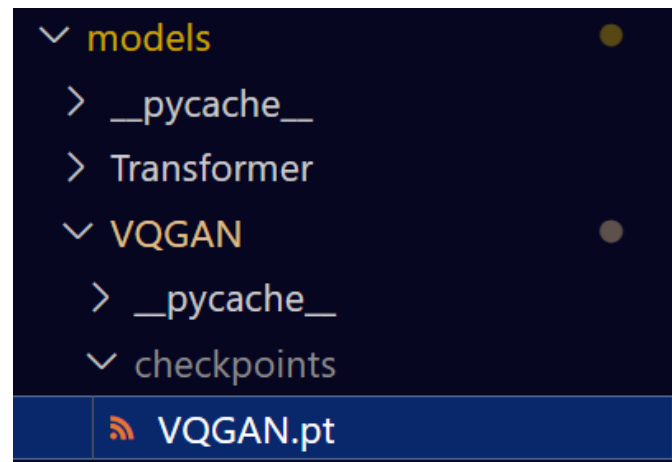- **Download**

  i. ON your own machine

  ?> sftp -P 10046 pp037@140.113.215.196 (passwd: pp037OnClass)

  ?> get VQGAN.pt

  ii. ON Provided machine

  ?> sftp pp037@192.168.201.46 (passwd: pp037OnClass)

  ?> get VQGAN.pt

```
∨ models
  > __pycache__
  > Transformer
  ∨ VQGAN
    > __pycache__
    ∨ checkpoints
      ⅏ VQGAN.pt
```

```
∨ VQGAN
  > __pycache__
  > checkpoints
  ∨ config
    ! Discriminator.yml
    ! VQGAN.yml
```

```yaml
---
MODEL_NAME: VQ_GAN

model_param:
  image_channels: 3
  enc_channels: [128, 128, 128, 256, 256, 512]
  dec_channels: [128, 128, 256, 256, 512]
  latent_dim: 256
  img_resolution: 64
  latent_resolution: 16
  num_codebook_vectors: 1024
  beta: 0.25
```
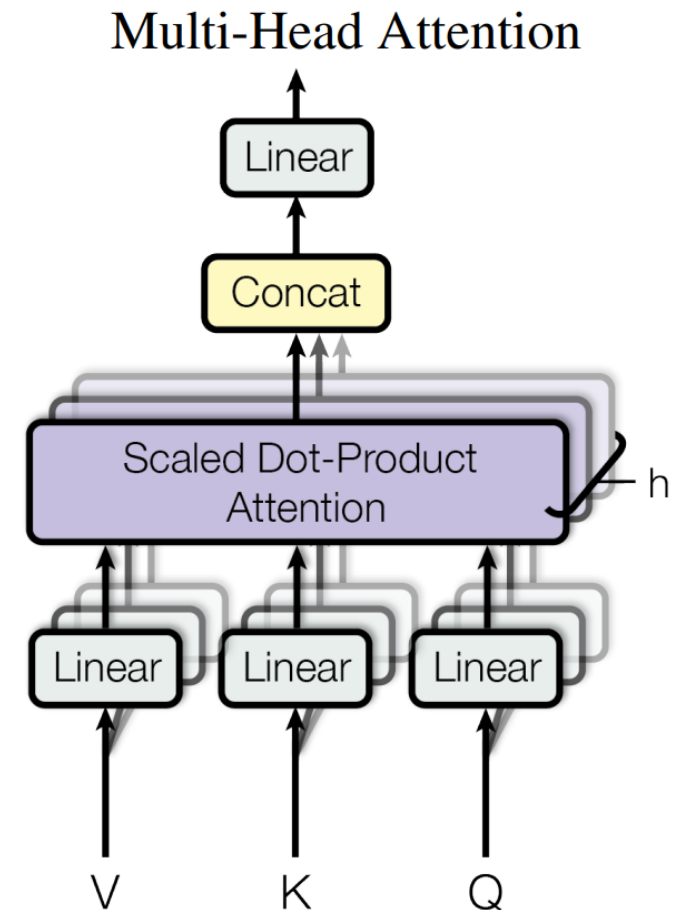
# Multi-Head Self-Attention

- You can't use any functions directly ex. torch.nn.MutiheadAttention

- Multi-Head Attention: total #s of head set to 16.

- Total $d_k$, $d_v$ set to 768

- $d_k$, $d_v$ for one head will be 768//16.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$



Multi-Head Attention

# MaskGIT Stage2 Training

- **You can't modify any model structure.**
- **Multi-Head Attention: total #s of head set to 16.**
-

```
∨ config
  ! MaskGit.yml
```

```yaml
---
MODEL_NAME: MaskGit


model_param:

  VQ_Configs:
    VQ_config_path: models/VQGAN/config/VQGAN.yml
    VQ_CKPT_path: models/VQGAN/checkpoints/VQGAN.pt


  num_image_tokens: 256
  num_codebook_vectors: 1024
  choice_temperature: 4.5
  gamma_type: cosine


  Transformer_param:
    num_image_tokens: 256
    num_codebook_vectors: 1024

    dim: 768
    n_layers: 15
    hidden_dim: 1536
```

- **How to set the Masked token?**

```python
class BidirectionalTransformer(nn.Module):
    def __init__(self, configs):
        super(BidirectionalTransformer, self).__init__()
        self.num_image_tokens = configs['num_image_tokens']
        #mask_token_id:1024
        self.tok_emb = nn.Embedding(configs['num_codebook_vectors'] + 1, configs['dim'])
        self.pos_emb = nn.init.trunc_normal_(nn.Parameter(torch.zeros(configs['num_image_tokens'], configs['dim'])), 0., 0.02)

        self.blocks = nn.Sequential(*[Encoder(configs['dim'], configs['hidden_dim']) for _ in range(configs['n_layers'])])
        self.Token_Prediction = TokenPredictor(configs['dim'])
        self.LN = nn.LayerNorm(configs['dim'], eps=1e-12)
        self.drop = nn.Dropout(p=0.1)

        self.bias = nn.Parameter(torch.zeros(self.num_image_tokens, configs['num_codebook_vectors'] + 1))
        self.apply(weights_init)

    def forward(self, x):
        # Token domain -> Latent domain
        token_embeddings = self.tok_emb(x)

        embed = self.drop(self.LN(token_embeddings + self.pos_emb))
        embed = self.blocks(embed)
        embed = self.Token_Prediction(embed)

        # Latent domain -> Token domain
        logits = torch.matmul(embed, self.tok_emb.weight.T) + self.bias

        return logits
```
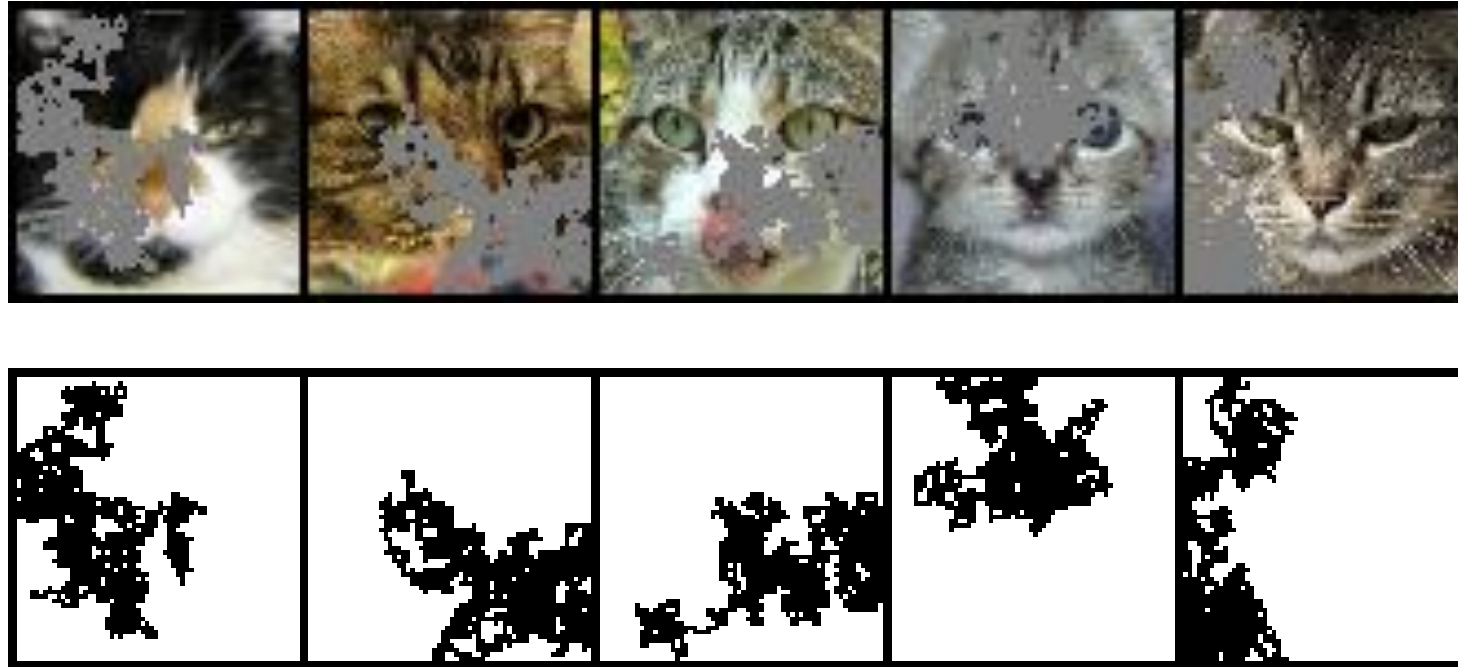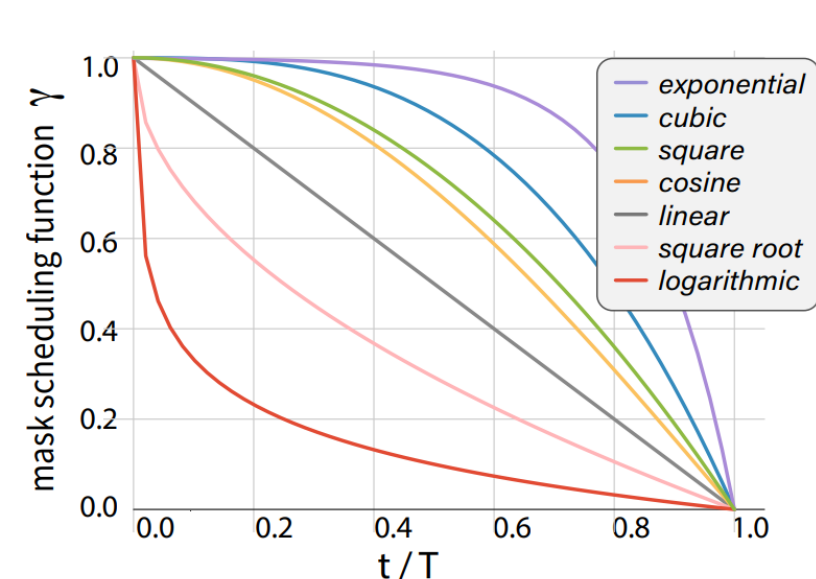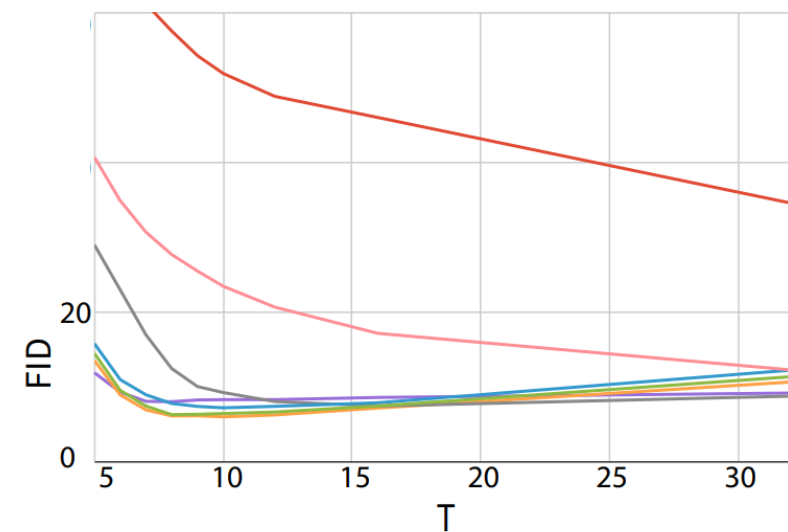
# Inference for Image Inpainting Task



- Tokenize the masked image
- Interpret the **inpainting mask** as the initial mask in iterative decoding

# Iterative Decoding



- **Mask Scheduling Functions** $\gamma(\frac{t}{T})$
  - •cosine    • linear    • square

- **Number of iterations** $T$
  (you can adjust)

- **Sweet spot** $t$
  (you can adjust)

| $\gamma$ | $T$ | FID ↓ | IS ↑ | NLL |
|---|---|---|---|---|
| Exponential | 8 | 7.89 | 156.3 | 4.83 |
| Cubic | 9 | 7.26 | 165.2 | 4.63 |
| Square | 10 | 6.35 | 179.9 | 4.38 |
| **Cosine** | 10 | **6.06** | **181.5** | 4.22 |
| Linear | 16 | 7.51 | 113.2 | 3.75 |
| Square Root | 32 | 12.33 | 99.0 | 3.34 |
| Logarithmic | 60 | 29.17 | 47.9 | 3.08 |

# Requirements

1. Download the dataset and pretrained weight of VQGAN (MaksGIT stage1).
2. Implement the Multi-head attention module on your own, if you use any function directly, your demo score will -10.
3. Train your transformer model (MaskGIT stage2) from scratch.
4. Implement iterative decoding for inpainting task.
5. Compare the FID score with different settings of mask scheduling parameters and visualize the iterative decoding for mask in latent domain or predicted images, if you don't show the visualization of iterative decoding when demo, your demo score will -20, meaning that you won't get any experiment score.

# Report Spec (40%)

**1. Introduction (5%)**

**2. Implementation Details (60%)**

    A. The details of your model (Multi-Head Self-Attention)

    B. The details of your stage2 training (MVTM, forward, loss)

    C. The details of your inference for inpainting task (iterative decoding)

**3. Experimental results (30%)**

    A. The best testing fid(21%)

        • Screenshot

        • Predicted image, Mask in latent domain with mask scheduling

        • The setting about training strategy, mask scheduling parameters, and so on

    B. Comparison figures with different mask scheduling parameters setting(total 9%) (each 3%)

        •cosine    • linear    • square

**4. Discussion(5%)**
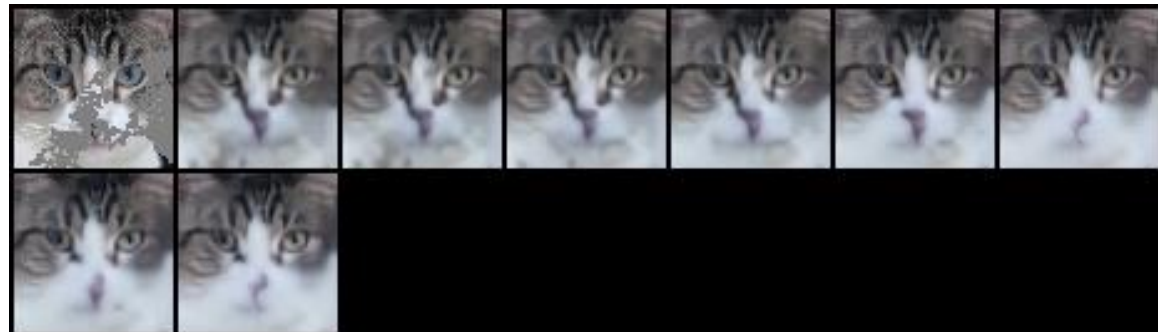
    A. Anything you want to share

# Demo (70%) (Prove your code implementation is correct)

- Show Multi-Head Attention module.
  - If you **directly use any functions**, your demo score will -10.

- Choose **either one** to show iterative decoding.
  - If **Both missing**, your demo score will -20.

**1.Mask in latent domain (specific 2 serial number)**



**2.Predicted image (specific 2 serial number)**

# Demo (70%)

Experiment Score

```
cd faster-pytorch-fid
python fid_score_gpu.py --predicted-path /path/your_inpainting_results_folder --device cuda:0
```

- **Experimental result (20%)**

| Average FID | Score |
|---|---|
| $40 \geq FID$ | 20 |
| $45 \geq FID > 40$ | 17 |
| $50 \geq FID > 45$ | 14 |
| $55 \geq FID > 50$ | 11 |
| $60 \geq FID > 55$ | 8 |
| $65 \geq FID > 60$ | 5 |
| $FID > 65$ | 0 |

- **Question (50%)**

# References

1. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In NeurIPS, 2017. https://arxiv.org/pdf/2202.04200.pdf

2. A. van den Oord, O. Vinyals, et al., "Neural discrete representation learning," in Advances in Neural Information Processing Systems, pp. 6306–6315, 2017. https://arxiv.org/abs/1711.00937

3. Esser, P., Rombach, R., and Ommer, B.: Taming Transformers for High-Resolution Image Synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12873–12883 (2021) https://arxiv.org/abs/2012.09841

4. Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T. Freeman. Maskgit: Masked generative image transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, June 2022. https://arxiv.org/abs/2202.04200