

Deep Learning Homework 3 Report

Cheng-Liang Chi

March 31, 2025

Contents

1	Introduction	2
2	Implementation Details	2
2.1	Multi-Head Self-Attention	2
2.2	Masked Visual Token Modeling (MVTM)	3
2.3	Iterative Decoding for Inpainting	4
3	Discussion	7
4	Experiment Score	7
4.1	Iterative Decoding	7
4.1.1	Mask in Latent Domain	7
4.1.2	Predicted Images	8
4.2	Best FID Score	8
4.2.1	Results	8
4.2.2	Experimental Settings	9

1 Introduction

In this lab, we explore the use of **MaskGIT** [1], a masked generative image transformer, for the **image inpainting** task. Image inpainting aims to restore missing or corrupted regions in an image in a visually plausible manner. To achieve this, we implement the second stage of the MaskGIT pipeline, which involves constructing a **Multi-Head Self-Attention** module [2], training a **Bidirectional Transformer** using **Masked Visual Token Modeling (MVTM)**, and performing **iterative decoding** to complete masked images.

The core idea behind MaskGIT is to leverage a bidirectional Transformer architecture to address the inefficiencies of traditional autoregressive models, enhancing both generation quality and computational efficiency. Initially, images are encoded using a pretrained VQGAN encoder [3] to generate discrete latent tokens. During training, random subsets of these tokens are masked, and the model learns to predict these masked tokens. During inference, iterative decoding is used, employing various mask scheduling strategies to progressively refine the reconstruction based on confidence scores.

This report details our implementation steps, key design choices, analyses of mask scheduling strategies, and evaluation of inpainting quality using the **FID score** [4]. Our experiments aim to illustrate MaskGIT's effectiveness in generating semantically consistent and visually high-quality image reconstructions.

2 Implementation Details

This section provides a detailed explanation of our implementation, which is structured into three parts: Multi-Head Self-Attention, Masked Visual Token Modeling (MVTM), and Iterative Decoding for Inpainting.

2.1 Multi-Head Self-Attention

In the Multi-Head Self-Attention implementation, the input tensor is first linearly projected into query (Q), key (K), and value (V) vectors. Each of these vectors is reshaped into multiple heads to allow parallel attention computation. We employ scaled dot-product attention independently on each head, which calculates attention scores by the dot product of Q and K, scaled by the square root of the dimension of the keys. After computing the softmax scores, they are applied to the values (V). The outputs from different heads are concatenated and passed through a final linear layer to produce the attention output.

Specific implementation details and dimension calculations can be referred to in the provided source code.

```
4 class MultiHeadAttention(nn.Module):
5     def __init__(self, dim=768, num_heads=16, attn_drop=0.1):
6         super(MultiHeadAttention, self).__init__()
7         self.head_num = num_heads
8         self.head_dim = dim // num_heads
9         self.scale = self.head_dim**-0.5
10
11         self.to_qkv = nn.Linear(dim, 3 * self.head_num * self.head_dim,
12                                   ↪ bias=False)
```

```

13     self.head = nn.Sequential(
14         nn.Linear(dim, dim),
15         nn.Dropout(attn_drop),
16     )
17
18     def forward(self, x):
19         """Hint: input x tensor shape is (batch_size, num_image_tokens,
20             ↳ dim),
21         because the bidirectional transformer first will embed each token
22             ↳ to dim dimension,
23         and then pass to n_layers of encoders consist of Multi-Head
24             ↳ Attention and MLP.
25         # of head set 16
26         Total d_k , d_v set to 768
27         d_k , d_v for one head will be 768//16.
28         """
29         qkv = self.to_qkv(x)
30         q, k, v = qkv.reshape(
31             3,
32             x.shape[0],
33             self.head_num,
34             x.shape[1],
35             self.head_dim,
36         )
37
38         attn = (q @ k.transpose(-2, -1)) * self.scale

```

2.2 Masked Visual Token Modeling (MVTM)

The MVTM module begins by encoding input images into latent tokens using a pretrained VQGAN encoder. We randomly mask a subset of these tokens based on a Bernoulli distribution, replacing masked tokens with a dedicated mask token identifier. The masked sequence is then passed through a bidirectional Transformer, which predicts the original tokens at masked positions. The training objective is to minimize the cross-entropy loss between the predicted tokens and the ground-truth latent tokens.

Specific details regarding mask token selection, masking ratios, and Transformer configurations are described explicitly in the source code implementation.

```

32     @torch.no_grad()
33     def encode_to_z(self, x):
34         codebook_mapping, codebook_indices, _ = self.vqgan.encode(x)
35         return codebook_mapping,
36             ↳ codebook_indices.reshape(codebook_mapping.shape[0], -1)

```

```

60     def forward(self, x):
61         _, z_indices = self.encode_to_z(x)
62         mask = torch.bernoulli(
63             0.5 * torch.ones(z_indices.shape, device=z_indices.device)
64         ).bool()
65         mask_token = (
66             self.mask_token_id
67             * torch.ones(z_indices.shape, device=z_indices.device).long()
68         )
69         masked_tokens = mask * mask_token + (~mask) * z_indices
70         logits = self.transformer(masked_tokens)
71         return logits, z_indices

```

2.3 Iterative Decoding for Inpainting

During inference, iterative decoding reconstructs masked regions of latent tokens through successive refinement. Initially, all masked positions are filled with the mask token. At each decoding iteration, the Transformer predicts token probabilities, and tokens with the highest confidence scores are fixed for subsequent iterations, while lower confidence tokens are remasked. Various mask scheduling strategies, including cosine, linear, and square, determine how the masking ratio decreases throughout the iterations. This iterative decoding continues until no masked tokens remain.

Implementation details such as temperature annealing, Gumbel noise application for sampling tokens, and specific mask scheduling functions are available within the provided implementation.

```

33     # mask_b: iteration decoding initial mask, where mask_b is true means
34     ↪ mask
35     def inpainting(self, image, mask_b, i): # MakGIT inference
36         maska = torch.zeros(
37             self.total_iter, 3, 16, 16
38         ) # save all iterations of masks in latent domain
39         imga = torch.zeros(
40             self.total_iter + 1, 3, 64, 64
41         ) # save all iterations of decoded images
42         mean = torch.tensor([0.4868, 0.4341, 0.3844],
43             ↪ device=self.device).view(3, 1, 1)
44         std = torch.tensor([0.2620, 0.2527, 0.2543],
45             ↪ device=self.device).view(3, 1, 1)
46         ori = (image[0] * std) + mean
47         imga[0] = ori # mask the first image be the ground truth of masked
48         ↪ image
49
50         self.model.eval()
51         with torch.no_grad():
52             _, z_indices = self.model.encode_to_z(image)
53             mask_num = mask_b.sum() # total number of mask token

```

```

50     z_indices_predict = z_indices
51     mask_bc = mask_b
52     mask_b = mask_b.to(device=self.device)
53     mask_bc = mask_bc.to(device=self.device)
54
55     ratio = 0
56     dec_img_ori = None
57     # iterative decoding for loop design
58     # Hint: it's better to save original mask and the updated mask
59     # ↳ by scheduling separately
60     for step in range(self.total_iter):
61         if step == self.sweet_spot:
62             break
63         ratio = step / self.total_iter
64
65         z_indices_predict, mask_bc = self.model.inpainting(
66             z_indices_predict,
67             mask_bc,
68             mask_num,
69             ratio,
70         )
71
72         # static method you can modify or not, make sure your
73         # ↳ visualization results are correct
74         mask_i = mask_bc.view(1, 16, 16)
75         mask_image = torch.ones(3, 16, 16)
76         indices = torch.nonzero(mask_i, as_tuple=False) # label
77         # ↳ mask true
78         mask_image[:, indices[:, 1], indices[:, 2]] = 0 # 3,16,16
79         maska[step] = mask_image
80         shape = (1, 16, 16, 256)
81         z_q =
82         # ↳ self.model.vqgan.codebook.embedding(z_indices_predict).view(shape)
83         z_q = z_q.permute(0, 3, 1, 2)
84         decoded_img = self.model.vqgan.decode(z_q)
85         dec_img_ori = (decoded_img[0] * std) + mean
86         imga[step + 1] = dec_img_ori # get decoded image
87
88     assert dec_img_ori is not None
89     ##decoded image of the sweet spot only, the test_results folder
90     # ↳ path will be the --predicted-path for fid score calculation
91     vutils.save_image(
92         dec_img_ori, os.path.join("test_results",
93         # ↳ f"image_{i:03d}.png"), nrow=1
94     )
95
96     # demo score
97     vutils.save_image(
98         maska,
99         os.path.join("mask_scheduling", f"test_{i}.png"),

```

```

71     return logits, z_indices
72
73 @torch.no_grad()
74 def inpainting(self, z_indices, mask, mask_num, step):
75     """predict"""
76     masked_z_indices = mask * self.mask_token_id + (~mask) * z_indices
77     logits = self.transformer(masked_z_indices)
78
79     # Apply softmax to convert logits into a probability distribution
80     ↪ across the last dimension.
81     prob = torch.softmax(logits, dim=-1)
82     # Find max probability for each token value & max prob index
83     z_indices_predict_prob, z_indices_predict = prob.max(dim=-1)
84
85     """ mask schedule """
86     ratio = self.gamma(step)
87     mask_len = torch.ceil(mask_num * ratio).long()
88
89     """ sample """
90     # predicted probabilities add temperature annealing gumbel noise as
91     ↪ confidence
92     g = torch.distributions.gumbel.Gumbel(0, 1).sample(
93         z_indices_predict_prob.shape
94     ) # gumbel noise
95     assert g is not None
96     g = g.to(z_indices_predict_prob.device)
97
98     temperature = self.choice_temperature * (1 - ratio)
99     confidence = z_indices_predict_prob + temperature * g
100
101     """ mask """
102     # hint: If mask is False, the probability should be set to
103     ↪ infinity, so that the tokens are not affected by the
104     ↪ transformer's prediction
105     confidence[~mask] = torch.inf
106     # sort the confidence for the rank
107     _, idx = confidence.topk(mask_len, dim=-1, largest=False)
108     # define how much the iteration remain predicted tokens by mask
109     ↪ scheduling
110     mask_bc = torch.zeros(
111         z_indices.shape, dtype=torch.bool,
112         device=z_indices_predict_prob.device
113     )
114     # At the end of the decoding process, add back the
115     ↪ original(non-masked) token values

```

3 Discussion

In this section, we discuss key insights gained from our experiments, potential limitations of our approach, and directions for future improvement.

My experiments showed the choice of mask scheduling function significantly influences the quality of inpainting results. The cosine scheduling consistently delivered the best performance, highlighting its effectiveness in gradually and smoothly reducing the mask ratio. Conversely, linear and square schedules showed slightly lower performance, possibly due to abrupt or uneven changes in masking ratios.

Another area worth exploring is the impact of different iterative decoding parameters, such as the number of iterations and the sweet spot. While my current settings provided good results, further experiments could optimize these parameters to balance computational efficiency and inpainting quality.

Lastly, future work might include applying our MaskGIT implementation to higher-resolution datasets or exploring other Transformer architectures to potentially improve inpainting quality and diversity further.

4 Experiment Score

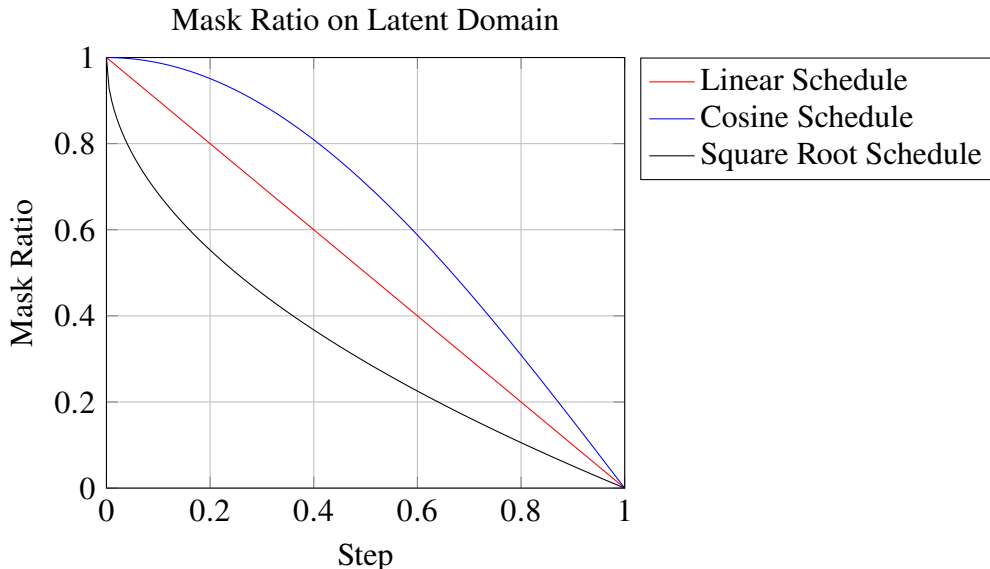
This section presents the experimental results validating our implementation, structured into two main parts: demonstration of iterative decoding and reporting of the best FID scores obtained.

4.1 Iterative Decoding

In this part, we demonstrate the iterative decoding process for different mask scheduling strategies.

4.1.1 Mask in Latent Domain

Present visualizations of the mask distribution over iterations using different scheduling strategies (cosine, linear, square).



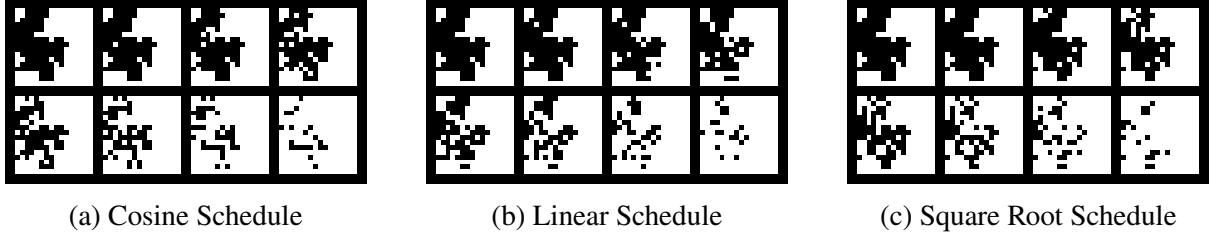


Figure 1: Mask in Latent Domain with Different Schedules. Three of this schedule are generated with the same image.

4.1.2 Predicted Images

The intermediate predicted images at different decoding iterations are shown in Figures 2a, 2b, and 2c. The images illustrate the model’s ability to progressively fill in the masked regions, demonstrating the effectiveness of the iterative decoding process.

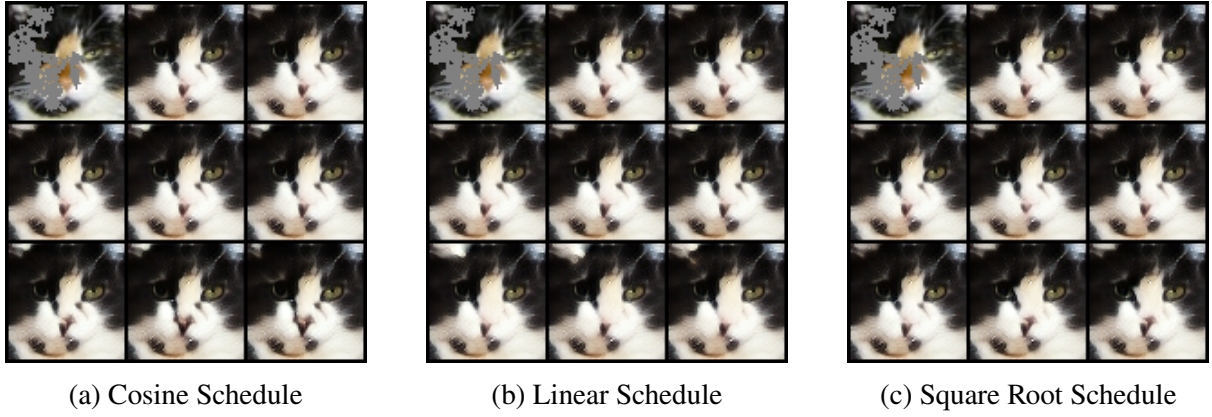


Figure 2: Predicted Images at Different Iterations. The images show the model’s iterative refinement process.

The predicted images for the linear schedule are shown in Figure 2b. The images demonstrate the iterative refinement process, where the model progressively fills in the masked regions.

4.2 Best FID Score

4.2.1 Results

The FID scores obtained from the experiments are summarized in Table 1. The scores indicate the quality of the generated images, with lower scores indicating better performance.

Table 1: FID Scores for Different Mask Scheduling Strategies

Mask Scheduling Strategy	FID Score
Cosine Schedule	38.02
Linear Schedule	39.12
Square Root Schedule	39.30

4.2.2 Experimental Settings

The experiments were conducted with the following settings:

- **Learning Rate:** 1×10^{-4}
- **Batch Size:** 32
- **Epochs:** 200
- **Sweet Spot:** 8
- **Total Iterations:** 8

And the results shown in Table 1 are using the lowest loss on the validation set.

References

- [1] H. Chang, H. Zhang, L. Jiang, C. Liu, and W. T. Freeman, “Maskgit: Masked generative image transformer,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 11 315–11 325. [Online]. Available: <https://arxiv.org/abs/2202.04200>.
- [2] A. Vaswani, N. Shazeer, N. Parmar, *et al.*, “Attention is all you need,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998–6008. [Online]. Available: <https://arxiv.org/abs/1706.03762>.
- [3] P. Esser, R. Rombach, and B. Ommer, “Taming transformers for high-resolution image synthesis,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 12 873–12 883. [Online]. Available: <https://arxiv.org/abs/2012.09841>.
- [4] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, “Gans trained by a two time-scale update rule converge to a local nash equilibrium,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 6626–6637. [Online]. Available: <https://arxiv.org/abs/1706.08500>.
- [5] Hank891008, *Deep-learning*, 2023. [Online]. Available: <https://github.com/hank891008/Deep-Learning>.