

# Supplementary Materials

Atsuya Nakata<sup>✉</sup> and Takao Yamanaka<sup>✉</sup>

Sophia University, Tokyo, Japan  
a-nakata-7r0@eagle.sophia.ac.jp, takao-y@sophia.ac.jp

## A Source codes of network architecture

The source codes for the network architecture in the proposed method are shown in Program 1.1 and Program 1.2. Program 1.1 is the source code of the low-resolution model at the first stage, where `MultiAxisTransformerLayer` is the MaxViT layer provided in Program 1.2. The input  $x$  is added with the positional encoding and the feature map of the conditional image compressed by CNN. The source code of the high-resolution model is almost same as Program 1.1, except that the feature map of the NFOV image extracted from the low-resolution image at the first stage is additionally added to the input  $x$ .

## B Comparison of models at first and second stages

The proposed method adopted MaxViT models both at the first and second stages for the low-resolution and high-resolution models. The different models were also examined for the models, as shown in Table 4. (1) is the proposed method using MaxViT at both stages. In (2), MaxViT at the first stage was replaced with Transformer. In (3), MaxViT at the second stage was replaced with MultiAxisTransformer [21]. In (4), Grid Attention in MaxViT at the second stage was replaced with the neighbor attention in Fig. 10. Although IS and LPIPS were highest with (4), FID was greatly improved with (1) where MaxViT was used at both stages. Therefore, the model (1) was adopted as the proposed method in this paper.

## C Sample images of proposed method compared with conventional method

Additional sample images of the synthesized omni-directional images in Fig. 7 are shown in Fig. 11. A comparison of synthesized omni-directional images between 2S-ODIS and 2S-ODIS 2day are shown in Fig. 12. Furthermore, additional sample images of the NFOV images toward the ground in Fig. 8 are shown in Fig. 13.

**Program 1.1:** Source code of MaxViT at first stage in PyTorch

```

1 class MultiAxisTransformer(nn.Module):
2     def __init__(self, vocab_size, d_model, layer=8, seq_len
3         =512):
4         # vocab_size: The number of VQGAN Codebooks
5         # d_model: The number of dimentions of model
6         # layer: The number of MaxViTLayers
7         # seq_len: The resolution of the input latent
8             variables
9         super(MultiAxisTransformerModel, self).__init__()
10        self.positional_encoding = nn.Parameter(torch.randn
11            (1, seq_len, d_model))
12        self.embedding = nn.Embedding(vocab_size+1, d_model)
13        self.trans_layers = nn.ModuleList([
14            MultiAxisTransformerLayer(d_model, 8) for i in
15            range(layer)])
16        self.mask_condition_conv = CNN(d_model)
17        self.fc = nn.Linear(d_model, vocab_size)
18
19    def forward(self, x, masked_condition, *args):
20        x = self.embedding(x)
21        condition = einops.rearrange(self.mask_condition_conv
22            (masked_condition), "b c h w->b (h w) c")
23        x = x + self.positional_encoding + condition
24        for layer in self.trans_layers:
25            x = layer(x)
26        output = self.fc(x)
27        return output

```

Program 1.2: Source code of MaxViT layer in PyTorch

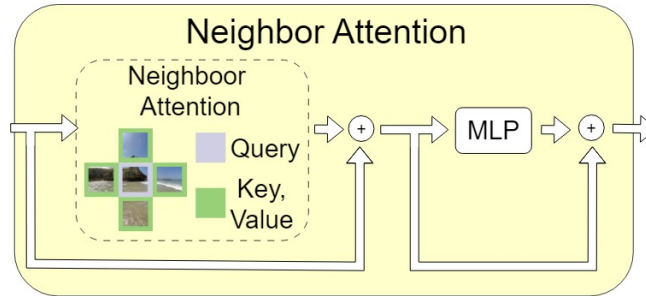
```

1 class MultiAxisTransformerLayer(nn.Module):
2     def __init__(self, d_model, n_head, patch_size=4):
3         # d_model: The number of dimentions of model
4         # n_head: Number of attention heads
5         # patch_size: The number of patch size of block
6         # attention
7
8         super().__init__()
9         self.mbconv = MBConvBlock(d_model)
10        self.block_attention = nn.TransformerEncoderLayer(
11            d_model, n_head, dim_feedforward=d_model*4,
12            batch_first=True, norm_first=True)
13        self.grid_attention = nn.TransformerEncoderLayer(
14            d_model, n_head, dim_feedforward=d_model*4,
15            batch_first=True, norm_first=True)
16        self.patch_size = patch_size
17
18    def forward(self, x):
19        x = einops.rearrange(x, "b (h w) c->b c h w", h=16)
20        b, c, h, w = x.shape
21        x = self.mbconv(x)
22        x = einops.rearrange(x, "b c (h1 h2) (w1 w2)->(b h1 w1
23            ) (h2 w2) c",
24            h1=h//self.patch_size, h2=self.patch_size,
25            w1=w//self.patch_size, w2=self.patch_size)
26        x = self.block_attention(x)
27        x = einops.rearrange(x, "(b h1 w1) (h2 w2) c->(b h2 w2
28            ) (h1 w1) c",
29            h1=h//self.patch_size, h2=self.patch_size,
30            w1=w//self.patch_size, w2=self.patch_size)
31        x = self.grid_attention(x)
32        x = einops.rearrange(x, "(b h2 w2) (h1 w1) c->b (h1 h2
33            w1 w2) c",
34            h1=h//self.patch_size, h2=self.patch_size,
35            w1=w//self.patch_size, w2=self.patch_size)
36        return x

```

**Table 4:** Comparison of models at first and second stages

First Stage	Second Stage	IS ( $\uparrow$ )	FID ( $\downarrow$ )	LPIPS ( $\uparrow$ )
(1) MaxViT	MaxViT	5.969	<b>18.263</b>	0.662
(2) Transformer	MaxViT	6.010	23.168	0.660
(3) MaxViT	MultiAxisTransformer [21]	5.820	22.274	0.657
(4) MaxViT	NeighborAttention (Fig. 10)	<b>6.030</b>	20.529	<b>0.667</b>

**Fig. 10:** Neighbor attention. The attention is calculated with adjacent N FoV images

## References

21. Zhao, L., Zhang, Z., Chen, T., Metaxas, D., Zhang, H.: Improved transformer for high-resolution gans. In: Advances in Neural Information Processing Systems (NeurIPS). vol. 34, pp. 18367–18380 (2021)

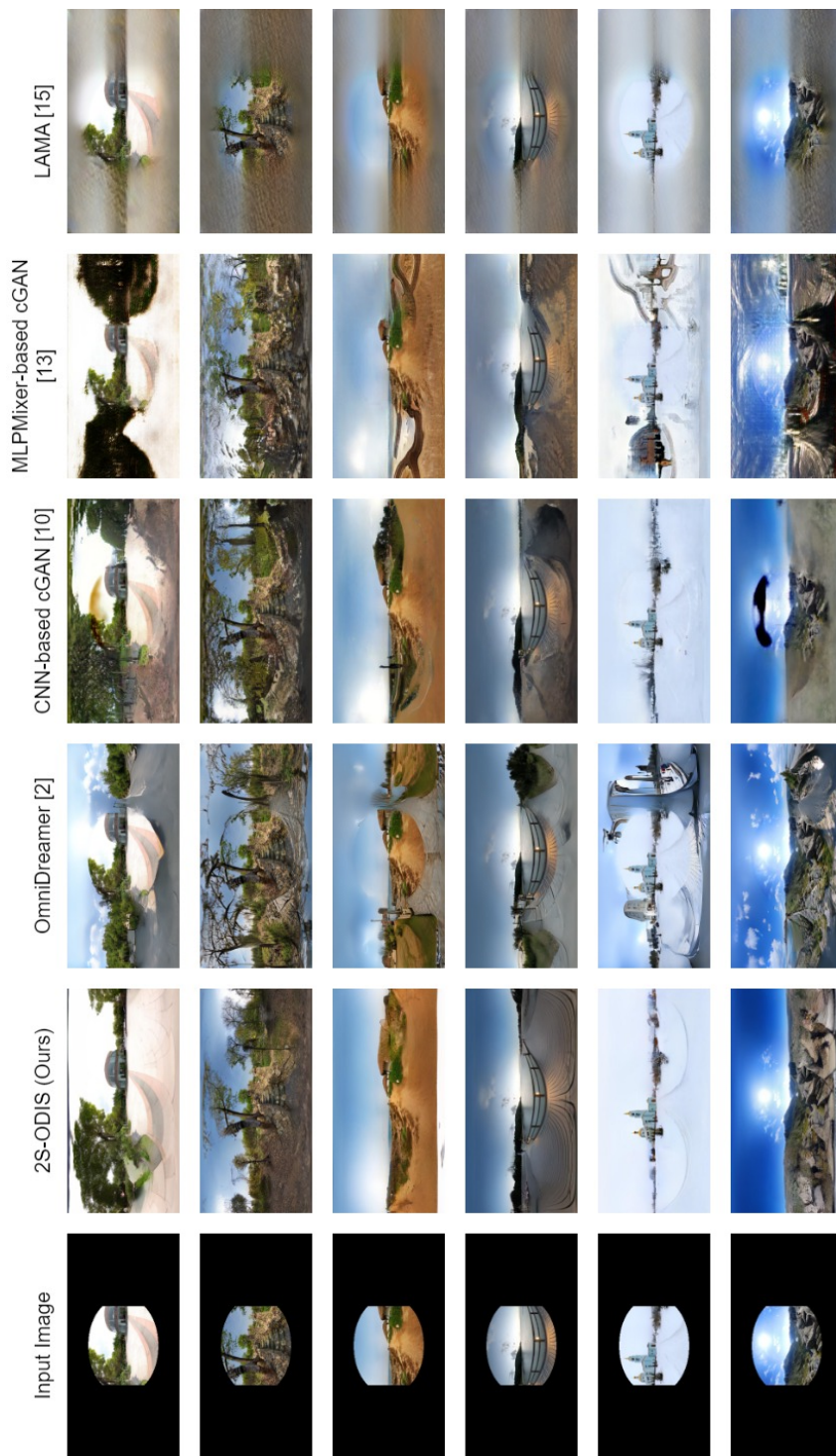
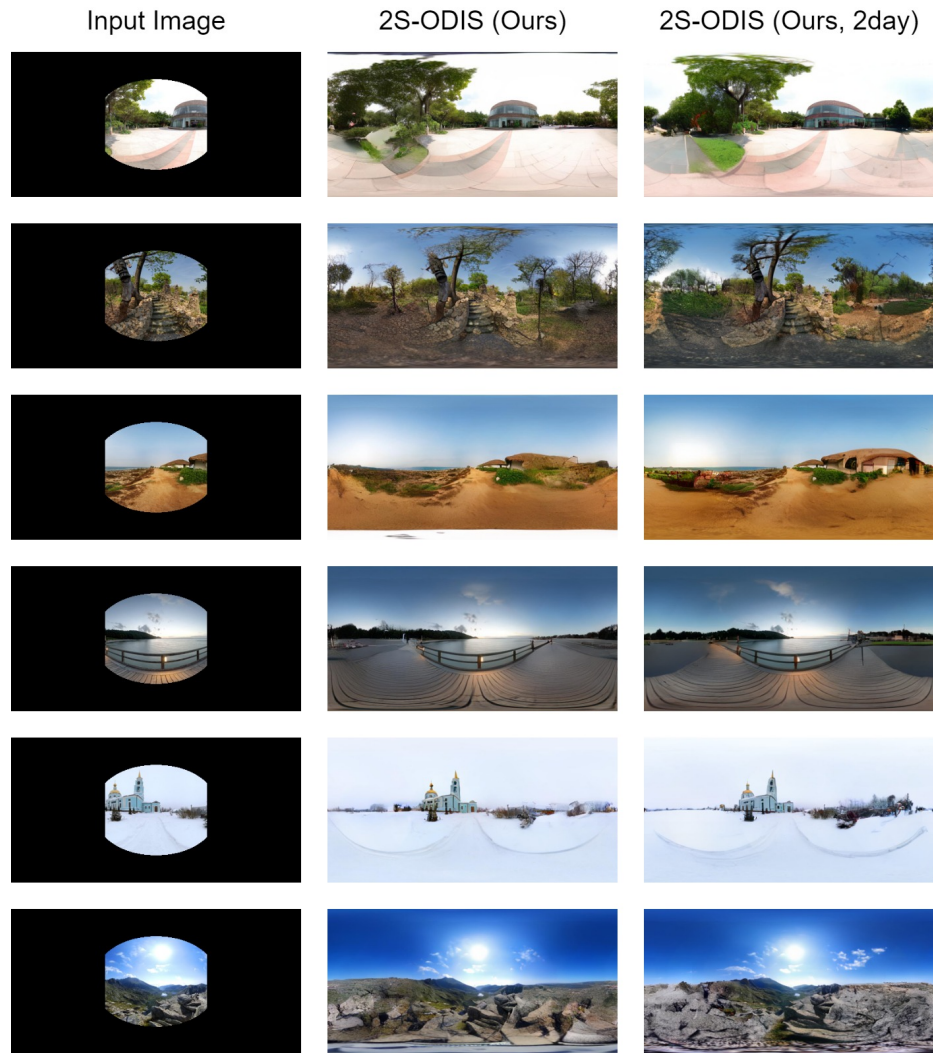
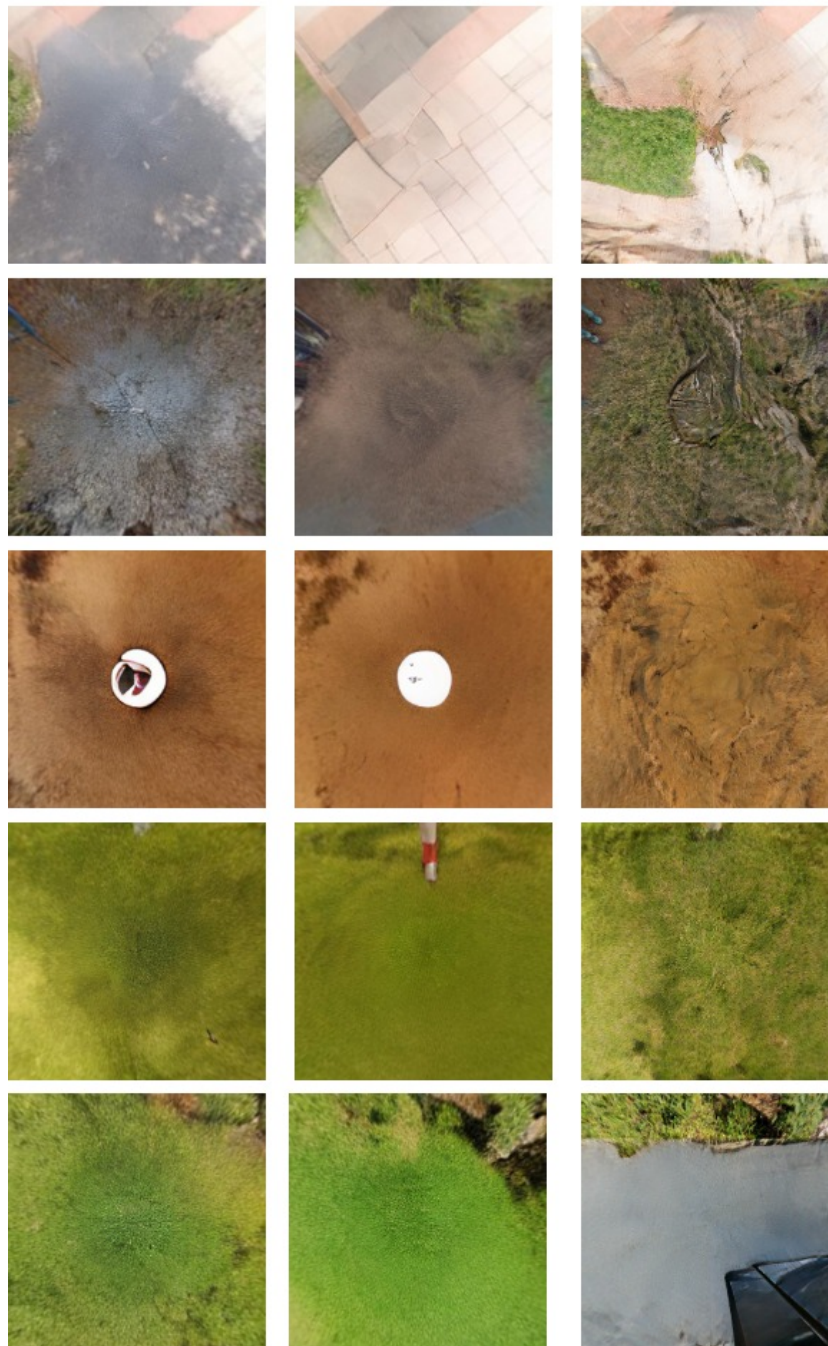


Fig. 11: Additional examples of synthesized omni-directional images in Fig. 7



**Fig. 12:** A comparison of synthesized omni-directional images between 2S-ODIS and 2S-ODIS 2day

2S-ODIS First Stage    2S-ODIS Second Stage    OmniDreamer [2]  
(Ours)



**Fig. 13:** Additional examples of NFOV images toward ground extracted from synthesized omni-directional images in Fig. 8.