# Appendix

## A. Visualization

### A.1. Qualitative results

Fig. A1 and Fig. A2 showcase qualitative comparisons with Atlantis on the D3 and D5 subsets of the Sea-thru [2] dataset and the SQUID [3] dataset. All models trained on the SynTIDE dataset, including AdaBins [4], NeWCRFs [8], PixerFormer [1], and MIM [7], consistently present better visual results on underwater images compared with those trained on the Atlantis dataset. Especially in the results of the first two close-shot images in Fig. A1, the model trained on the Atlantis dataset fails to clearly show the difference in distance between the ball and the background. In contrast, our results match the ground truth closer, more distinctly displaying the contrast between the ball and the background in the image.
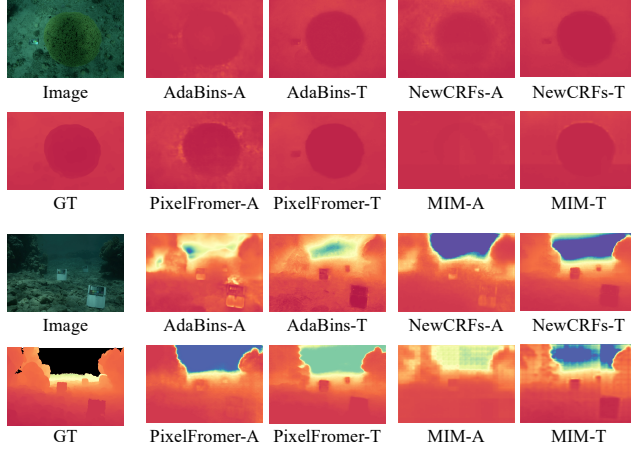


Figure A1. Qualitative results on the Sea-thru dataset [2]. '-A' and '-T' denote models trained on Atlantis [9] and Our SynTIDE dataset, respectively. The depth estimation results are notably improved after training on our dataset. Due to the original 'Image' being extremely dim, the content is hardly visible. To clearly display the content of 'Image', we adjust its contrast and brightness in this figure. These adjustments do not apply to any inference or evaluation processes at the code level.

### A.2. Zero-shot underwater depth data generation

Thanks to our training strategy, which fine-tunes the pretrained text-to-image model [5] using LoRA [6] with a minor low rank, we retain its strong generalization ability to a certain extent. This enables TIDE to generate underwater depth data for scenes and objects never seen during training, as shown in Fig. A3. Even when the provided text prompts contain objects that do not exist in the real world, such as Godzilla, TIDE can still generate seemingly reasonable underwater image-depth pairs. However, this ca-
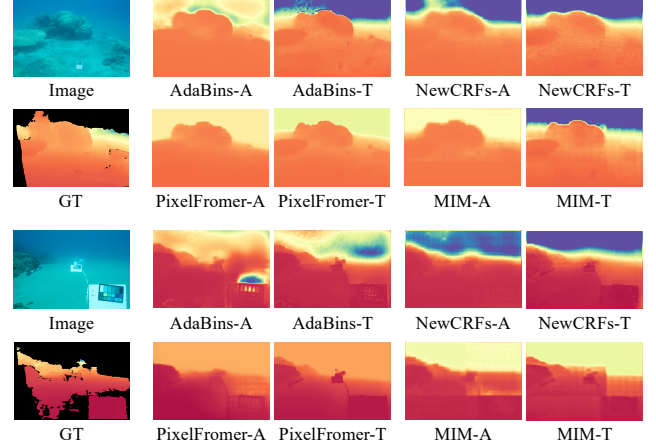


Figure A2. Qualitative results on SQUID dataset [3]. '-A' and '-T' denote models trained on Atlantis [9] and Our SynTIDE dataset, respectively. The depth estimation results are notably improved after training on our dataset.

pability is particularly challenging for Atlantis [9], which requires the depth map in advance as a condition.

## B. Pseudo-code

To demonstrate the simplicity of TIDE and each component, we provide pseudo code with PyTorch in Listing 1, Listing 2, and Listing 3. These codes are simple and easy to implement. The complete code to reproduce the experiments will be made available before the conference.

An underwater photo of a Godzilla.

An underwater photo of an avocado.

An underwater photo of a red crocodile.

An underwater photo of a table.

An underwater photo of
a huge underwater glacier.

An underwater photo of a penguin.

An underwater photo of a UFO.

Realistic, luxurious underwater hotel.
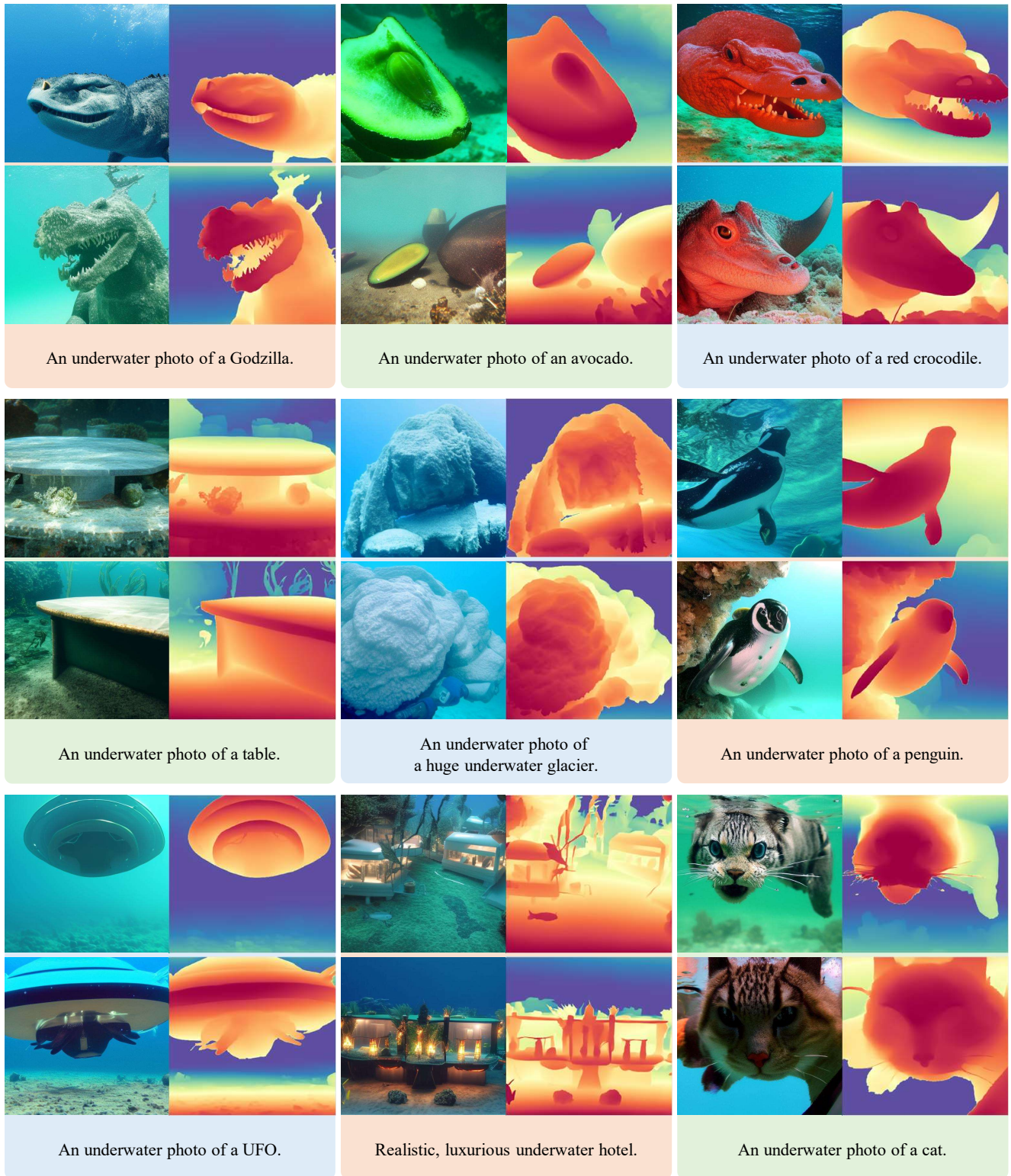
An underwater photo of a cat.

Figure A3. Representative zero-shot image-depth pairs synthesized by TIDE present strong consistency, diversity, and generalization. Images of relevant categories are not included in the training data.

```
1    from torch import nn
2
3    class TAN(nn.Module):
4        def __init__(self, nhidden=1152, hidden_dim=256, time_hidden_dim=6912):
5            super().__init__()
6
7            self.gamma_mlp = nn.MLP(input_dim=nhidden, hidden_dim=hidden_dim, output_dim=nhidden)
8            self.beta_mlp = nn.MLP(input_dim=nhidden, hidden_dim=hidden_dim, output_dim=nhidden)
9            self.time_adaptive_scale = nn.Sequential(
10               nn.Linear(time_hidden_dim, 1),
11               nn.Sigmoid(),
12           )
13
14       def forward(self, x, time_embed, modal1_feats, modal2_feats=None):
15           if modal2_feats is not None:
16               gamma1, beta1 = self._forward(modal1_feats)
17               gamma2, beta2 = self._forward(modal2_feats)
18
19               gamma = (gamma1 + gamma2) / 2
20               beta = (beta1 + beta2) / 2
21           else:
22               gamma, beta = self._forward(modal1_feats)
23
24           sigma = self.time_adaptive_scale(time_embed)
25           out = x * (1 + sigma * gamma) + sigma * beta
26           return out
27
28       def _forward(self, modal_feats):
29           gamma = self.gamma_mlp(modal_feats)
30           beta = self.beta_mlp(modal_feats)
31           return gamma, beta
```

Listing 1. TAN PyTorch code. When multiple modal features are input, they share the MLP weights and average the multiple sets of gamma and beta.

```
1    from torch import einsum
2    import torch.nn.functional as F
3
4    def ILS_Attention(attn, image_feats, text_feats, cross_attn_map=None):
5        query = attn.to_q(image_feats)
6        key = attn.to_k(text_feats)
7        value = attn.to_v(text_feats)
8
9        if cross_attn_map is not None:
10           hidden_states = einsum('b h l n, b h n c -> b h l c', cross_attn_map, value)
11       else:
12           hidden_states, cross_attn_map = F.scaled_dot_product_attention(query, key, value)
13
14       return hidden_states, cross_attn_map
```

Listing 2. The ILS Attention PyTorch code. ILS_Attention is a component within the Transformer/MiniTransformer block in Listing 3.

```
1   from torch import nn
2   from pixart import PixartTransfomer, Transformer, MiniTransformer
3
4   class TIDE(PixartTransfomer):
5       def __init__(self, transfomer_layer=28, mini_transformer_layer=10):
6           super().__init__()
7           self.t2i_transformer = Transformer(num_layer=transfomer_layer)
8           self.t2d_transformer = MiniTransformer(num_layer=mini_transformer_layer)
9           self.t2m_transformer = MiniTransformer(num_layer=mini_transformer_layer)
10
11          self.D2M_tan_blocks = nn.ModuleList(
12              [
13                  TAN(nhidden=1152, hidden_dim=256)
14                  for _ in range(mini_transformer_layer)
15              ]
16          )
17          self.M2D_tan_blocks = nn.ModuleList(
18              [
19                  TAN(nhidden=1152, hidden_dim=256)
20                  for _ in range(mini_transformer_layer)
21              ]
22          )
23          self.DM2I_tan_blocks = nn.ModuleList(
24              [
25                  TAN(nhidden=1152, hidden_dim=256)
26                  for _ in range(mini_transformer_layer)
27              ]
28          )
29
30          self.dense_blocks_inject_pos = [0, 3, 6, 9, 12, 15, 18, 21, 24, 27]
31
32      def forward(self, hidden_states_image, hidden_states_depth, hidden_states_mask,
33              hidden_states_text, time_embed
34      ):
35          for block_index, t2i_block in enumerate(self.t2i_transformer.blocks):
36              hidden_states_image, implicit_layout = t2i_block(
37                  hidden_states_image,
38                  encoder_hidden_states=hidden_states_text,
39                  time_embed=time_embed,
40              )
41              if block_index in self.dense_blocks_inject_pos:
42                  id = self.dense_blocks_inject_pos.index(block_index)
43
44                  hidden_states_mask = self.D2M_tan_blocks[id](
45                      hidden_states_mask, time_embed, hidden_states_depth
46                  )
47
48                  hidden_states_depth, _ = self.t2d_transformer.blocks[id](
49                      hidden_states_depth,
50                      encoder_hidden_states=hidden_states_text,
51                      time_embed=time_embed,
52                      implicit_layout=implicit_layout,
53                  )
54
55                  hidden_states_mask, _ = self.t2m_transformer.blocks[id](
56                      hidden_states_mask,
57                      encoder_hidden_states=hidden_states_text,
58                      time_embed=time_embed,
59                      implicit_layout=implicit_layout,
60                  )
61
62                  hidden_states_depth = self.M2D_tan_blocks[id](
63                      hidden_states_depth, time_embed, hidden_states_mask
64                  )
65
66                  hidden_states_image = self.DM2I_tan_blocks[id](
67                      hidden_states_image, time_embed, hidden_states_depth, hidden_states_mask
68                  )
69
70          image_noise_output, depth_noise_output, mask_noise_output = self.output(
71              hidden_states_image, hidden_states_depth, hidden_states_mask
72          )
73          return image_noise_output, depth_noise_output, mask_noise_output
74
```

Listing 3. TIDE PyTorch code. The entire code of TIDE will be made available before the conference.

# References

[1] Ashutosh Agarwal and Chetan Arora. Attention attention everywhere: Monocular depth prediction with skip attention. In *Proc. of IEEE Winter Conf. on Applications of Computer Vision*, pages 5861–5870, 2023. 1

[2] Derya Akkaynak and Tali Treibitz. Sea-thru: A method for removing water from underwater images. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 1682–1691, 2019. 1

[3] Dana Berman, Deborah Levy, Shai Avidan, and Tali Treibitz. Underwater single image color restoration using hazelines and a new quantitative dataset. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(8):2822–2837, 2020. 1

[4] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 4009–4018, 2021. 1

[5] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-$\alpha$: Fast training of diffusion transformer for photorealistic text-to-image synthesis, 2024. 1

[6] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *Proc. of Intl. Conf. on Learning Representations*, 2022. 1

[7] Zhenda Xie, Zigang Geng, Jingcheng Hu, Zheng Zhang, Han Hu, and Yue Cao. Revealing the dark secrets of masked image modeling. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 14475–14485, 2023. 1

[8] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. Neural window fully-connected crfs for monocular depth estimation. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 3916–3925, 2022. 1

[9] Fan Zhang, Shaodi You, Yu Li, and Ying Fu. Atlantis: Enabling underwater depth estimation with stable diffusion. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 11852–11861, 2024. 1