

A Unified Image-Dense Annotation Generation Model for Underwater Scenes

Hongkai Lin Dingkang Liang Zhenghao Qi Xiang Bai*
 Huazhong University of Science and Technology
 {hklin, dkliang, xbai}@hust.edu.cn

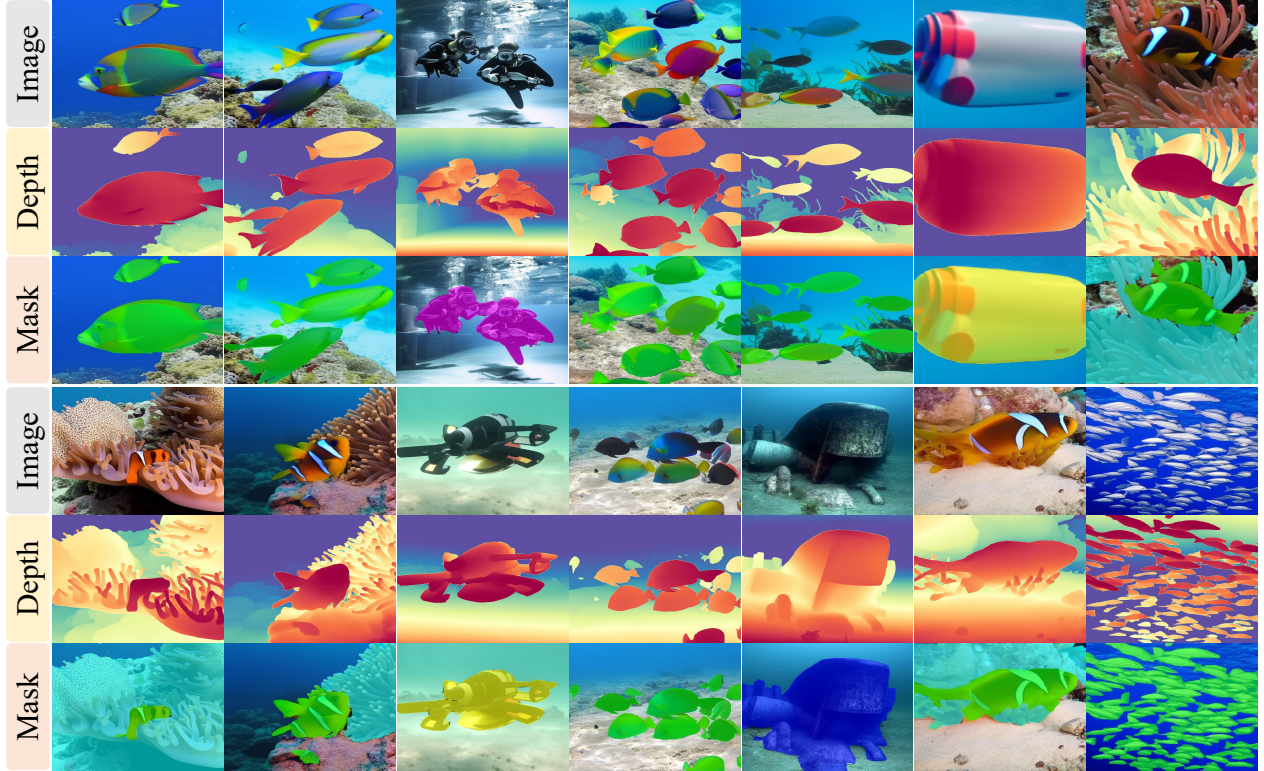


Figure 1. **We present TIDE, a unified underwater image-dense annotation generation model.** Its core lies in the shared layout information and the natural complementarity between multimodal features. Our model, derived from the text-to-image model and fine-tuned with underwater data, enables the generation of highly consistent underwater image-dense annotations from solely text conditions.

Abstract

Underwater dense prediction, especially depth estimation and semantic segmentation, is crucial for gaining a comprehensive understanding of underwater scenes. Nevertheless, high-quality and large-scale underwater datasets with dense annotations remain scarce because of the complex environment and the exorbitant data collection costs. This paper proposes a unified **Text-to-Image** and **Dense** annotation generation method (TIDE) for underwater scenes. It relies solely on text as input to simultaneously generate realistic underwater images and multiple highly consistent dense annotations. Specifically, we unify the generation of

text-to-image and text-to-dense annotations within a single model. The **Implicit Layout Sharing mechanism (ILS)** and cross-modal interaction method called **Time Adaptive Normalization (TAN)** are introduced to jointly optimize the consistency between image and dense annotations. We synthesize a large-scale underwater dataset using TIDE to validate the effectiveness of our method in underwater dense prediction tasks. The results demonstrate that our method effectively improves the performance of existing underwater dense prediction models and mitigates the scarcity of underwater data with dense annotations. We hope our method can offer new perspectives on alleviating data scarcity issues in other fields. The code is available at <https://github.com/HongkaiLin/TIDE>.

* Corresponding author.

1. Introduction

Underwater dense prediction, particularly depth estimation and semantic segmentation, is essential for underwater exploration and environmental monitoring. However, the complex environment and the prohibitive data collection costs result in a scarcity of underwater data with dense annotations. Such conditions severely hinder the advancement of dense prediction technologies in underwater scenes.

Fortunately, the recent success of the image generative technique [14, 29, 43] provides a breakthrough in addressing the scarcity of underwater scene data. In the field of general object understanding, controllable data synthesis [17, 23, 31, 38] demonstrates its effectiveness in few-shot scenarios. A straightforward solution is to apply them to underwater scenes directly. For instance, Atlantis [42], a pioneering controllable generation method for underwater depth data that takes ControlNet as its core, utilizes terrestrial depth maps as conditions. It effectively mitigates the issue of scarce underwater depth data and achieves consistent performance improvements across multiple underwater depth datasets and models.

Despite remarkable progress, there are still challenges in Atlantis, as follows: 1) Atlantis, as shown in Fig. 2(a), generates underwater depth data using terrestrial depth maps as conditions due to the lack of underwater depth maps. It is considered a suboptimal approach since it may not align with natural underwater scenes. Better recreating authentic underwater environments is equally essential. 2) It generates data with only a single type of dense annotations, which is insufficient for understanding complex underwater scenes. Thus, a natural question arises: *How can we simultaneously generate highly consistent, one-to-many, and vivid underwater images and dense annotation pairs?*

In this paper, we explore the possibility of simultaneously generating highly consistent, realistic underwater scene images and multiple types of dense annotations using only text conditions. Our approach, which we refer to as **TIDE**, is illustrated in Fig. 2(b), presents a unified **Text-to-Image** and **Dense** annotation generation method. TIDE is an end-to-end training and inference model that integrates denoising models in parallel for both text-to-image generation and text-to-dense annotation generation.

To align the images and multiple type dense annotations generated by parallel denoising models, we propose the **Implicit Layout Sharing (ILS)** mechanism. Specifically, the cross-attention map as an implicit layout is the key to controlling the image layout in the text-to-image model [6, 29], inspiring us to share the implicit layout for aligning images and dense annotations. ILS effortlessly replaces the cross-attention map in the text-to-dense annotation model with that from the text-to-image model, effectively improving the consistency between the image and dense annotations. Furthermore, considering the intrinsic

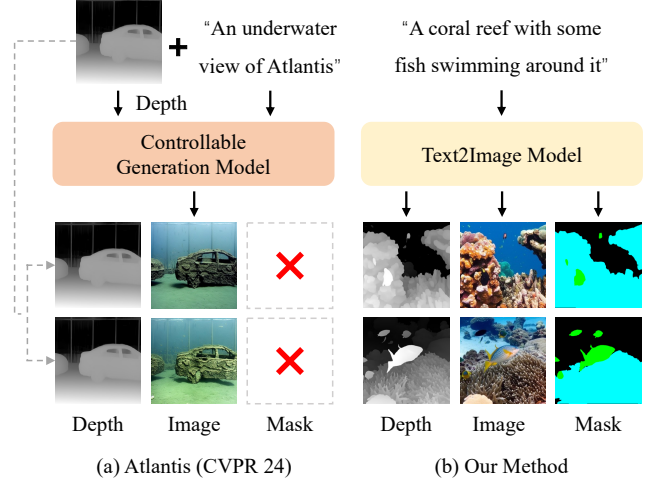


Figure 2. The comparison between Atlantis [42] and our method. Unlike Atlantis, which requires text and depth map conditions, our method only needs text as the input condition to generate image-dense annotations (e.g., depth maps and semantic masks).

complementarity between features of different modalities, we introduce a cross-modal interaction method called **Time Adaptive Normalization (TAN)**, a normalization layer that modulates the activations using different modal features. The consistency of the image and dense annotations can further be jointly optimized through cross-modal feature interaction among different dense annotation generation and between image and dense annotation generation.

To verify the effectiveness of our method, we use TIDE to generate a large-scale dataset of underwater images with dense annotations named SynTIDE. Extensive experiments demonstrate the effectiveness of SynTIDE for underwater dense prediction tasks. In the underwater depth estimation task, SynTIDE presents consistent improvements in various fine-tuning models. For example, when adopting representative NewCRFs [41] as the fine-tuning model, our approach achieves significance gains over previous work, particularly in the SI_{log} and δ_1 metrics, with improvements of 14.73 and 36% on the D3 and D5 subsets of Sea-thru [3] dataset, respectively. In underwater semantic segmentation, pre-training with SynTIDE yields consistent improvements across different models. For instance, when using ViT-Adapter [7] as the training model, pre-training with the SynTIDE dataset leads to improvements of 2.1% mIoU on the USIS10K [21] dataset.

TIDE demonstrates powerful data generation capabilities for underwater scenes. Using only easily accessible text prompts, TIDE can generate highly consistent and realistic underwater images and multiple types of dense annotations. It holds potential as a mainstream data synthesis method for underwater scenes and offers a promising direction for alleviating data scarcity in other fields. The main contribu-

tions of this work are as follows: **1)** We propose a novel data synthesis method, TIDE, which uses text as the sole condition to generate images and their corresponding multi-type dense annotations simultaneously. To our knowledge, TIDE is the first method capable of simultaneously synthesizing both images and multiple dense annotations from text. **2)** To align the images and dense annotations, we introduce the Implicit Layout Sharing mechanism. The text-to-image and text-to-dense annotation models share the same layout information, ensuring proper alignment between the image and dense annotations. Meanwhile, the consistency between image and dense annotations can be further optimized through the cross-modal interaction method called Time Adaptive Normalization.

2. Relate Work

2.1. Underwater Dense Prediction

Dense prediction tasks in underwater scenes are crucial for comprehensively understanding underwater scenes. The publication of SUIM [16] provides a fundamental dataset and benchmark for the exploration of underwater semantic segmentation. To fill the gap in underwater instance segmentation, WaterMask [20] publishes the UIIS dataset, and a model is designed to cater to the unique characteristics of underwater images, improving the accuracy of underwater instance segmentation. Recently, the rise of general foundational segmentation models [18, 28] drives further development in the field of underwater segmentation [21, 37, 44].

Due to the lack of underwater depth estimation datasets, most underwater depth estimation methods focus on traditional techniques, unsupervised, or self-supervised approaches. Traditional methods [9] mainly rely on statistical priors, such as the dark channel prior [12], to estimate underwater depth. Gupta et al. [10] model the relationship between underwater and above-water hazy appearances to depth estimation. UW-GAN [11] and Atlantis [42] improve the performance of underwater depth estimation by synthesizing training datasets through generative models.

While these methods make notable contributions to underwater dense prediction tasks, the large-scale and high-quality dataset in underwater scenes with only segmentation or depth annotations remains insufficient for achieving comprehensive underwater scene understanding.

2.2. Controllable Data Synthesis

Thanks to the success of diffusion models [14] and the availability of large-scale, high-quality text-image training data, text-to-image models [6, 24, 27, 29] and controllable image generation models [22, 36, 43] achieve unprecedented success in image quality, diversity, and consistency.

He et al. [13] are the first to explore and demonstrate the effectiveness of state-of-the-art text-to-image genera-

tion models for image recognition. This makes it possible to achieve diverse data collection and accurate annotation at a lower cost. Wu et al. and Nguyen et al. [23, 32] explore the ability of pre-trained diffusion models to enhance real data in few-shot settings for segmentation tasks. Diffu-mask [33] ingeniously combines text-to-image models with AffinityNet [2], achieving open-vocabulary segmentation data synthesis. Freemask [38] demonstrates that synthetic data can further enhance the performance of semantic segmentation models under fully supervised settings by incorporating freestyle, a controllable image generation method using semantic masks as input conditions. Seggen [40] designs a multi-stage semantic segmentation data synthesis method, text2mask and mask2image, which achieves high semantic consistency semantic segmentation data only using text as the condition. Detdiffusion [31] synthesizes object detection data by incorporating object categories and spatial coordinates into the text.

Unlike the aforementioned single-task data synthesis methods, we propose a novel end-to-end underwater data synthesis approach that simultaneously generates semantic masks and depth maps, relying solely on text conditions.

3. Preliminaries

Diffusion Models (DMs) [14] emerge as leading text-to-image (T2I) generation models, recognized for their ability to produce realistic images. DMs can reconstruct data distribution by learning the reverse process of a diffusion process. Denoting z_t as the random variable at t -th timestep, the diffusion process is modeled as a Markov Chain:

$$z_t \sim \mathcal{N}(\sqrt{\alpha_t}z_{t-1}, (1 - \alpha_t)\mathbf{I}), \quad (1)$$

where α_t is the fixed coefficient predefined in the noise schedule, and \mathbf{I} refers to identity matrix. A prominent variant, the Latent Diffusion Model (LDM) [29], innovatively shifts the diffusion process of standard DMs into a latent space. This transition notably decreases computational costs while preserving the generative quality and flexibility of the original model. The resulting efficiency gain primarily arises from the reduced dimensionality of the latent space, which allows for lower training costs without compromising the model’s generative capabilities.

Stable Diffusion, an exemplary implementation of LDM, comprises an AutoEncoder [30] and a latent diffusion model. The AutoEncoder \mathcal{E} is designed to learn a latent space that is perceptually equivalent to the image space. Meanwhile, the LDM ϵ_θ is parameterized as a denoising model with cross-attention and trained on a large-scale dataset of text-image pairs via:

$$\mathcal{L}_{LDM} := \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim N(0,1), t} [\|\epsilon - \epsilon_\theta(z_t, t, \tau_\theta(y))\|_2^2], \quad (2)$$

where ϵ is the target noise. τ_θ and y are the pre-trained

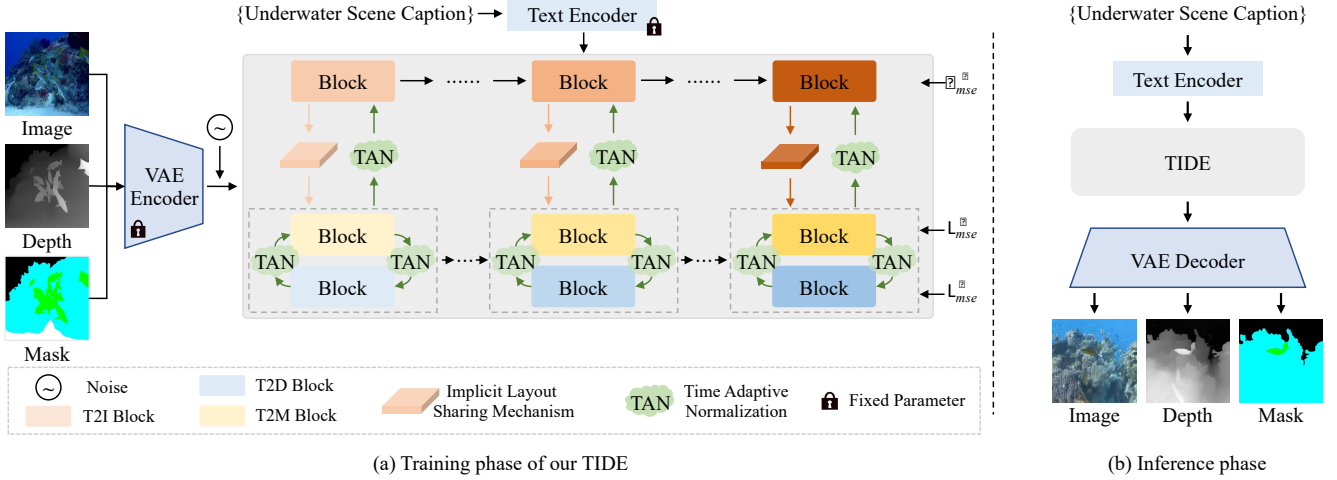


Figure 3. Training and Inference. The denoising model of TIDE mainly consists of three transformers, each dedicated to text-to-image, text-to-depth, and text-to-mask. The proposed Implicit Layout Sharing mechanism (ILS) and Time Adaptive Normalization (TAN) are used to align the generated image, depth map, and semantic mask.

text encoder (e.g., CLIP [25], T5 [26]) and text prompts, respectively. This equation represents the mean-squared error (MSE) between the target noise ϵ and the noise predicted by the model, encapsulating the core learning mechanism of the latent diffusion model.

4. Our Method

An overview of our method, a unified text-to-image and dense annotation generation model (TIDE), is shown in Fig. 3. TIDE is built upon a pre-trained transformer [6] for text-to-image generation, along with two fine-tuned mini-transformers (details provided in Sec. 5.2) dedicated to text-to-depth and text-to-mask generation. Simply parallelizing multiple text-to-image processes does not ensure consistency between the images and dense annotations. To enable consistency between them, we propose Implicit Layout Sharing (ILS) and the cross-modal interaction method named Time Adaptive Normalization (TAN). After training, TIDE simultaneously generates images and multiple dense annotations with high consistency using only text as input.

4.1. Data Preparation

We aim to generate realistic underwater images, corresponding highly consistent depth maps, and semantic masks. However, existing high-quality, dense annotation data primarily consists of mask annotations. Therefore, we construct training data around these datasets with semantic masks, as shown in Tab. 1. On this basis, we obtain the corresponding depth map and caption for each image using existing foundation models. Specifically, for each underwater image, the corresponding depth map is obtained by pre-trained Depth Anything [39]. Meanwhile, the caption of

each image is obtained from the pre-trained BLIP2 [19]. We construct approximately 14K quadruples {Image, Depth, Mask, Caption} for TIDE training.

Table 1. Segmentation Datasets and Data Splits. * denotes the training set of TIDE, while the others are used for evaluation.

Datasets	Seg Task	Train	Val	Test
SUIM [16]	Semantic	1,488*	110*	/
UIIS [20]	Instance	3,937*	691	/
USIS10K [21]	Instance	7,442*	1,594	1,596*

4.2. Implicit Layout Sharing Mechanism

In advanced text-to-image models [6, 29], the cross-attention map plays a crucial role in controlling the image layout. Existing methods [22, 36] demonstrate that adjusting the cross-attention map during the text-to-image process can effectively control the layout of the generated image. Therefore, the cross-attention map can be considered as the implicit layout information. Intuitively, sharing the implicit layout between text-to-image and text-to-dense annotations may establish a strong correlation between the generated image and dense annotations. To this end, we propose an Implicit Layout Sharing mechanism to align the generated image and dense annotations. Specifically, cross-attention, as a crucial process for generating implicit layouts in text-to-image/mask/depth model, can first be formulated as:

$$\begin{aligned}
 \text{Attn}_i(Q_i, K_i, V_i) &= \text{softmax}(Q_i K_i^\top / \sqrt{c}) V_i, \\
 \text{Attn}_d(Q_d, K_d, V_d) &= \text{softmax}(Q_d K_d^\top / \sqrt{c}) V_d, \\
 \text{Attn}_m(Q_m, K_m, V_m) &= \text{softmax}(Q_m K_m^\top / \sqrt{c}) V_m,
 \end{aligned} \quad (3)$$

where c refers to the feature channel. $Q_i/Q_d/Q_m$, $K_i/K_d/K_m$, and $V_i/V_d/V_m$ represent the query, key, and value within the text-to-image/depth/mask cross-attention module, respectively. Since text-to-image models are pre-trained on high-quality and large-scale image-caption datasets, they exhibit strong controllability and generalization. Therefore, sharing the implicit layouts from the text-to-image model is the optimal choice to ensure the quality of the generated data. As shown in Fig. 3(a), the implicit layouts from the block in the text-to-image model are shared with the cross-attention in the block of text-to-dense annotation models. The implicit layouts refer to:

$$M_i = \text{softmax}(Q_i K_i^T / \sqrt{c}). \quad (4)$$

By sharing the implicit layouts from the text-to-image model, the cross-attention of text-to-depth (Attn_d) and text-to-mask (Attn_m) can be simplified as follows:

$$\begin{aligned} \text{Attn}_d(Q_d, K_d, V_d) &= M_i \times V_d, \\ \text{Attn}_m(Q_m, K_m, V_m) &= M_i \times V_m, \end{aligned} \quad (5)$$

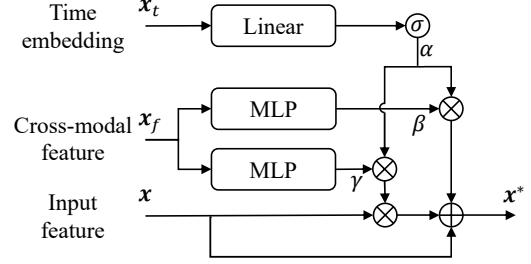
where \times refers to matrix multiplication. Implicit Layout Sharing is an elegant and efficient method that unifies image and dense annotation generation, improving consistency between them. It also reduces the overall generation cost, as there is no need to compute separate cross-attention maps for the text-to-dense annotation models.

4.3. Time Adaptive Normalization

Considering the complementary nature of different modality features, we propose a cross-modal feature interaction method called Time Adaptive Normalization (TAN), as shown in Fig. 4.

Specifically, TAN is utilized to adjust the image layout by leveraging the cross-modal features x_f from different branches. The cross-modal features are mapped to two normalization parameters, γ and β , by MLPs, which are used to control the variation in the image layout. In this context, the features from text-to-depth and text-to-mask serve as cross-modal input features for each other. For instance, in the TAN corresponding to the i -th text-to-depth block, the outputs from both the i -th text-to-depth block and the i -th text-to-mask block serve as the input feature x and cross-modal input feature x_f , respectively. A slight difference is that for text-to-image, the features from both text-to-depth and text-to-mask serve as the cross-modal features. In the TAN cross-modal interaction process of text-to-image, two sets of γ and β are obtained, provided by the different modalities features from text-to-depth and text-to-mask. These two sets of parameters are averaged to the $\bar{\gamma}$ and $\bar{\beta}$. Then, time embeddings x_t is introduced to adaptively control the influence of the cross-modal features. The normalization can be formalized as follows:

$$x' = \alpha \cdot \gamma x + \alpha \cdot \beta, \quad x^* = x' + x, \quad (6)$$



⊙ Sigmoid ⊕ Element-wise add ⊗ Element-wise product

Figure 4. In TAN, the cross-modal features are first mapped to the modulation parameters γ and β . Then, a time-adaptive confidence α is introduced to control the degree of normalization.

where x , x' , and x^* are the input feature, normalized feature, and output feature, respectively. α is the time adaptive coefficients obtained from x_t through linear transformation and sigmoid. The TAN will be applied not only from text-to-dense annotations to text-to-image but also between text-to-depth and text-to-mask to improve the consistency among dense annotations. Implicit Layout Sharing and Time Adaptive Normalization are two complementary methods that construct a joint interaction process, optimizing the consistency between the generated image and dense annotations during training.

4.4. Learning Objective

During training, the learnable parameters include only the proposed TAN module and the LoRA [15] used to fine-tune the pre-trained transformer. The overall loss \mathcal{L} is composed equally of the denoising losses from the three branches: text-to-image, text-to-depth, and text-to-mask:

$$\mathcal{L} = \mathcal{L}_{mse}^I + \mathcal{L}_{mse}^D + \mathcal{L}_{mse}^M. \quad (7)$$

4.5. Data Synthesis

Thanks to the proposed ILS and TAN, TIDE can generate realistic and highly consistent underwater images and dense annotations after training, using only text conditions, as shown in Fig. 3(b).

We filter out redundant parts from the 14K captions obtained in Sec. 4.1, resulting in approximately 5K non-redundant captions as text conditions. For each caption, we generate ten samples to construct a large-scale synthetic dataset named SynTIDE. Some representative examples are shown in Fig. 1. The SynTIDE dataset is utilized to validate the effectiveness of our method in the dense prediction task for underwater scenes.

4.6. Analysis

Insights of framework design. In the text-to-image model, the cross-attention map contains the layout information of

Table 2. Quantitative comparisons on real underwater depth estimation datasets.

Method	Fine-tuning Dataset	Reference	$SI_{log} \downarrow$	$A.Rel \downarrow$	$log_{10} \downarrow$	$RMSE \downarrow$	$S.Rel \downarrow$	$RMSE_{log} \downarrow$	$\delta_1 \uparrow$	$\delta_2 \uparrow$	$\delta_3 \uparrow$
Quantitative comparisons on the D3 and D5 subsets of Sea-thru [3] dataset.											
AdaBins [5]	Atlantis [42]	CVPR 24	38.24	1.33	0.12	1.41	12.89	0.39	0.50	0.81	0.92
	SynTIDE (Ours)	-	26.92 _(-11.32)	1.31 _(-0.02)	0.08 _(-0.04)	1.12 _(-0.29)	15.74 _(+2.85)	0.27 _(-0.12)	0.71 _(+0.21)	0.95 _(+0.14)	0.99 _(+0.07)
NewCRFs [41]	Atlantis [42]	CVPR 24	37.10	1.68	0.12	1.44	14.76	0.38	0.48	0.84	0.95
	SynTIDE (Ours)	-	22.37 _(-14.73)	1.50 _(-0.18)	0.06 _(-0.06)	1.24 _(-0.20)	22.50 _(+7.74)	0.23 _(-0.15)	0.84 _(+0.36)	0.97 _(+0.13)	0.99 _(+0.04)
PixelFormer [1]	Atlantis [42]	CVPR 24	23.70	1.34	0.06	1.17	17.29	0.24	0.81	0.97	0.99
	SynTIDE (Ours)	-	21.39 _(-2.31)	1.46 _(+0.12)	0.05 _(-0.01)	1.15 _(-0.02)	21.79 _(+4.50)	0.22 _(-0.02)	0.88 _(+0.07)	0.98 _(+0.01)	0.99 _(+0.00)
MIM [35]	Atlantis [42]	CVPR 24	37.01	1.37	0.11	1.51	14.42	0.38	0.56	0.84	0.94
	SynTIDE (Ours)	-	22.49 _(-14.52)	1.27 _(-0.10)	0.06 _(-0.05)	1.01 _(-0.50)	16.46 _(+2.04)	0.23 _(-0.15)	0.85 _(+0.29)	0.97 _(+0.13)	0.99 _(+0.05)
Quantitative comparisons on the SQUID [4] dataset.											
AdaBins [5]	Atlantis [42]	CVPR 24	29.56	0.28	0.11	2.24	0.69	0.31	0.56	0.86	0.94
	SynTIDE (Ours)	-	25.63 _(-3.93)	0.23 _(-0.05)	0.09 _(-0.02)	2.69 _(+0.45)	0.92 _(+0.23)	0.27 _(-0.04)	0.67 _(+0.11)	0.90 _(+0.04)	0.97 _(+0.03)
NewCRFs [41]	Atlantis [42]	CVPR 24	25.19	0.23	0.09	2.56	0.83	0.26	0.68	0.90	0.96
	SynTIDE (Ours)	-	25.55 _(+0.36)	0.23 _(-0.00)	0.09 _(+0.00)	3.02 _(+0.46)	1.07 _(+0.24)	0.27 _(+0.01)	0.68 _(+0.00)	0.91 _(+0.01)	0.97 _(+0.01)
PixelFormer [1]	Atlantis [42]	CVPR 24	21.34	0.18	0.07	1.86	0.43	0.22	0.76	0.94	0.98
	SynTIDE (Ours)	-	19.08 _(-2.26)	0.16 _(-0.02)	0.07 _(-0.00)	1.75 _(-0.11)	0.36 _(-0.07)	0.19 _(-0.03)	0.79 _(+0.03)	0.97 _(+0.03)	0.99 _(+0.01)
MIM [35]	Atlantis [42]	CVPR 24	27.45	0.26	0.10	2.14	0.68	0.28	0.61	0.88	0.95
	SynTIDE (Ours)	-	26.98 _(-0.47)	0.25 _(-0.01)	0.09 _(-0.01)	3.04 _(+0.90)	1.11 _(+0.43)	0.28 _(-0.00)	0.65 _(+0.04)	0.89 _(+0.01)	0.96 _(+0.01)

the image. Thus, the cross-attention map can be viewed as an implicit layout. If two text-to-image models share the same implicit layout and undergo proper fine-tuning, the generated images are likely to exhibit strong layout similarity. Therefore, we share the same implicit layout across multiple text-to-image models. Meanwhile, we use LoRA to fine-tune the multiple text-to-image models [6].

Zero-shot generation ability. Thanks to our training strategy, which fine-tunes the pre-trained text-to-image model using only LoRA, the generalization ability of the text-to-image model is retained to some extent. This enables TIDE to generate underwater images during inference that are not seen during training. Furthermore, due to the proposed Implicit Layout Sharing and Time Adaptive Normalization mechanisms, the generated depth maps align well with these images. Therefore, TIDE has the ability to generate zero-shot underwater image-depth map pairs.

5. Experiments

5.1. Dataset and Evaluation Metrics

Underwater Depth Estimation. We follow the work [42], the D3 and D5 subsets of Sea-thru [3], and the SQUID dataset [4] used to evaluate the depth estimation capability in underwater scenes. These datasets include underwater images with depth maps obtained via the Structure-from-Motion (SfM) algorithm.

The quantitative evaluation metrics include root

mean square error ($RMSE$) and its logarithmic variant ($RMSE_{log}$), absolute error in log-scale (log_{10}), absolute relative error ($A.Rel$), squared relative error ($S.Rel$), the percentage of inlier pixels (δ_i) with thresholds of 1.25^i , and scale-invariant error in log-scale (SI_{log}): $100\sqrt{Var(\epsilon_{log})}$.

Underwater Semantic Segmentation. The UIIS [20] and USIS10K [21] datasets are chosen to validate the effectiveness of our method in underwater semantic segmentation tasks. Instance masks belonging to the same semantic category are merged to construct semantic segmentation annotations for the UIIS and USIS10K datasets.

We calculate the mean Intersection over Union (mIoU) for six categories (i.e., Fish, Reefs, Aquatic Plants, Wrecks, Human Divers, and Robots) to evaluate the accuracy of the segmentation results.

5.2. Implementation Details

The training process consists of two parts: pre-training the mini-transformer and training TIDE. In the first stage, the mini-transformer is initialized with the first ten layers of the PixArt- α [6] pre-trained transformer. Then, the mini-transformer is trained for 60K iterations on the text-to-image task with all parameters. The training data consists of 14K underwater image-caption pairs from Sec. 4.1. In the second stage, the PixArt- α pre-trained transformer and the mini-transformer are used as initial weights for the text-to-image and text-to-dense annotation models, respectively.

Table 3. Quantitative results of underwater semantic segmentation.

Method	Backbone	Training Data		mIoU	
		Real	SynTIDE	UIIS	USIS10K
Segformer [34] (NeurIPS 21)	MiT-B4	✓		70.2	74.6
			✓	76.5	72.8
Mask2former [8] (CVPR 22)	Swin-B	✓		72.7	76.1
			✓	74.2	72.9
ViT-Adapter [7] (ICLR 23)	Adapter-B	✓		73.5	74.6
			✓	75.7	72.6
		✓	✓	75.1(+1.6)	76.7(+2.1)

Meanwhile, they are fine-tuned using LoRA [15] for 200K iterations with a batch size of 4. The LoRA ranks of the text-to-image/depth/mask branches are 32, 64, and 64, respectively. All experiments are conducted on a server with four NVIDIA 4090 24G GPUs.

5.3. Main results

Underwater Depth Estimation. We train four representative depth estimation models, Adasbin [5], NewCRFs [30], PixelFormer [1], and MIM [35], to present quantitative results, as shown in Tab. 2. Compared to previous underwater data synthesis work Atlantis [42], depth estimation models trained on our SynTIDE dataset show consistent improvements across most quantitative metrics on two evaluated datasets. Especially on MIM [35], a powerful pre-trained model, our method reduces the SI_{log} metric from $37.01 \rightarrow 22.49$ (-14.52) and improves δ_1 from $0.56 \rightarrow 0.85$ (+0.29) on the D3 and D5 subsets of the Sea-thru dataset. Meanwhile, on PixelFormer [1], a depth estimation model with outstanding generalization that also performs best for Atlantis, our method achieves better performance across nearly all quantitative metrics on both evaluated underwater depth estimation datasets.

These results demonstrate that our method achieves highly competitive consistency compared to Atlantis, which uses stronger dense conditions. Furthermore, the data generated by TIDE is closer to natural underwater scenes and shows rich species diversity. Most importantly, TIDE unifies the generation of images and multiple highly consistent dense annotations, capabilities that Atlantis lacks.

Underwater Semantic Segmentation. In the underwater semantic segmentation task, we validate the effectiveness of our method by pre-training with the SynTIDE dataset in three representative semantic segmentation models, Seg-

Table 4. Ablation on the impact of each TIDE component.

ILS	TAN	$SI_{log} \downarrow$	$A.Rel \downarrow$	$\delta_1 \uparrow$	mIoU
✓		24.46	1.23	0.76	36.8
	✓	24.59	1.40	0.78	36.2
✓	✓	23.71	1.37	0.79	42.1

Table 5. Ablation on the impact of the component positions. “{Start, End}” indicates the starting and ending positions where the operations are applied, with a step size of 3.

{Start, End}	$SI_{log} \downarrow$	$A.Rel \downarrow$	$\delta_1 \uparrow$	mIoU
{0,12}	22.06	1.46	0.86	34.7
{15,27}	44.86	1.09	0.43	8.6
{0,27}	23.71	1.37	0.79	42.1

Table 6. Ablation on the impact of scaling synthetic data for underwater dense prediction tasks.

N Sample	$SI_{log} \downarrow$	$A.Rel \downarrow$	$\delta_1 \uparrow$	mIoU
1	23.49	1.46	0.82	55.3
3	22.94	1.54	0.85	60.9
6	22.96	1.56	0.85	63.6
10	22.37	1.50	0.84	64.2

former [34], Mask2former [8], and ViT-Adapter [7]. Following the work [38], we filter the noise in the generated annotations with 1.5 tolerance.

Pre-training on high-quality synthetic datasets is widely recognized as a way to gain strong prior knowledge. On the UIIS dataset, models trained on the SynTIDE dataset consistently achieve superior results compared to real data. On the other larger USIS10K dataset, by further fine-tuning the model on the UIIS10K train set, we achieve notable improvements. Especially on ViT-Adapter, we enhance the performance of the model from $74.6\% \rightarrow 76.7\%$.

These results show that models pre-trained on the SynTIDE dataset exhibit strong prior knowledge in the underwater semantic segmentation task. Additionally, these results demonstrate that the unified image and dense annotation generation model proposed in this paper can generate highly consistent image-dense annotation pairs, making it suitable for various underwater dense prediction tasks.

5.4. Ablation Studies

Unless otherwise specified, we conduct ablation studies by training TIDE for 30K iterations. We synthesize three samples for each caption, as described in Sec. 4.5. We conduct ablation studies on the USIS10K dataset with SegFormer-B4 for semantic segmentation and the D3 and D5 subsets of

the Sea-thru dataset with NewCRFs for depth estimation.

Ablation on the effectiveness of each component. We first evaluate the contribution of each component within TIDE, as shown in Tab. 4. When utilizing only the Implicit Layout Sharing (ILS) mechanism or Time Adaptive Normalization (TAN), the former outperforms the latter in depth estimation and semantic segmentation. Combining both methods results in a significant improvement (36.8% \rightarrow 42.1%) in semantic segmentation. These results indicate that ILS and TAN are complementary methods. By combining them for end-to-end training, the consistency between images and dense annotations can be further optimized. Additionally, we further demonstrate the effectiveness of the time-adaptive operation. As shown in the last row of Tab. 4, without time-adaptive parameters, the quality of the generated data will be varying degrees of degradation, especially for the semantic segmentation task.

Ablation on the position of components. We then study the effect of the position of ILS and TAN, as shown in Tab. 5. We find that applying the ILS and TAN mechanisms in the first half of the transformer of text-to-image yields better performance than using them in the second half. This can be attributed to the layout information produced in the first half of the transformer, which is mismatched with the ILS introduced in the latter part. Meanwhile, the results demonstrate that combining both achieves better consistency between the image and dense annotations.

Ablation on data scaling. Finally, we synthesize N samples for each caption to validate the impact of synthetic data scale on underwater dense prediction tasks, as shown in Tab. 6. It can be observed that as the amount of synthetic data increases, there is no substantial improvement in the underwater depth estimation task. However, for the underwater semantic segmentation task, a significant gain is observed in the early stages as N increases, but the tendency of improvement begins to flatten after $N = 6$.

5.5. More Challenging Underwater Data Synthesis

We validate whether TIDE can generate more challenging data by adding extra text prompts about underwater lighting or water quality (e.g., low light, turbidity) to the original underwater scene caption. As shown in Fig. 5, the results demonstrate that TIDE can generate more challenging underwater images. While annotating these underwater images may be extremely difficult for humans, TIDE can effortlessly produce highly consistent and accurate dense annotations, which hold great practical value for real-world underwater applications. In addition, to demonstrate the diversity of generated underwater data, we generate twelve underwater images from the same text prompt, as shown in Fig. 6. It can be observed that, despite sharing the same text prompt, the generated images exhibit rich diversity.

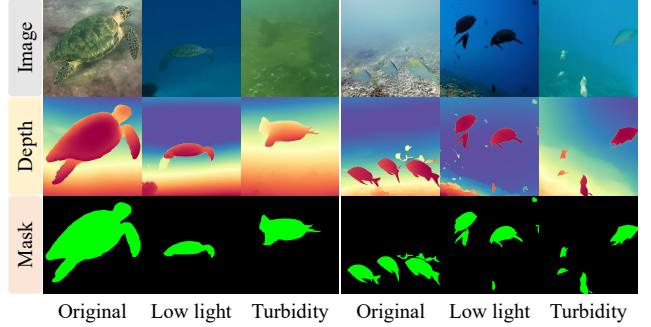


Figure 5. More challenging underwater data generated by TIDE.



Figure 6. Visualization of generated data diversity.

5.6. Limitation

Despite the promising results achieved, our method still has some limitations. First, our approach cannot generate instance-level semantic masks from the generation perspective. Relying on text prompts to guide the generation of instance-level masks with semantic annotations remains challenging. Additionally, although TIDE can leverage the powerful priors of pre-trained T2I models to generate highly challenging underwater images (e.g., low light, turbidity), there is still room for improvement. These will be key directions for future expansion.

6. Conclusion

This paper introduces a unified text-to-image and dense annotation generation model for underwater scenes. The model can generate realistic underwater images and multiple highly consistent dense annotations using only text prompts as input. We validate the effectiveness of our method on underwater depth estimation and semantic segmentation tasks by synthesizing a large-scale underwater dataset. In the depth estimation task, extensive experiments show that our method, using only text as input, achieves highly competitive results compared to previous methods that required stronger dense conditions for underwater depth synthesis. Meanwhile, pre-training with data synthesized using our method further improves model performance in the semantic segmentation task. Our study provides a new perspective for alleviating data scarcity.

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. 6, 7, 11
- [2] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 4981–4990, 2018. 3
- [3] 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. 2, 6, 11
- [4] Dana Berman, Deborah Levy, Shai Avidan, and Tali Treibitz. Underwater single image color restoration using haze-lines and a new quantitative dataset. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(8):2822–2837, 2020. 6, 11
- [5] 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. 6, 7, 11
- [6] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis, 2024. 2, 3, 4, 6, 11
- [7] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. In *Proc. of Intl. Conf. on Learning Representations*, 2023. 2, 7
- [8] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022. 7
- [9] Paulo LJ Drews, Erickson R Nascimento, Silvia SC Botelho, and Mario Fernando Montenegro Campos. Underwater depth estimation and image restoration based on single images. *IEEE computer graphics and applications*, 36(2):24–35, 2016. 3
- [10] Honey Gupta and Kaushik Mitra. Unsupervised single image underwater depth estimation. In *Proc. of IEEE Intl. Conf. on Image Processing*, pages 624–628. IEEE, 2019. 3
- [11] Praful Hambarde, Subrahmanyam Murala, and Abhinav Dhall. Uw-gan: Single-image depth estimation and image enhancement for underwater images. *IEEE Transactions on Instrumentation and Measurement*, 70:1–12, 2021. 3
- [12] Kaiming He, Jian Sun, and Xiaoou Tang. Single image haze removal using dark channel prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(12):2341–2353, 2010. 3
- [13] Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan Qi. Is synthetic data from generative models ready for image recognition? In *Proc. of Intl. Conf. on Learning Representations*, 2023. 3
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Proc. of Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2, 3
- [15] 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. 5, 7, 11
- [16] Md Jahidul Islam, Chelsey Edge, Yuyang Xiao, Peigen Luo, Muntaqim Mehtaz, Christopher Morse, Sadman Sakib Enan, and Junaed Sattar. Semantic segmentation of underwater imagery: Dataset and benchmark. In *Proc. of the IEEE Int. Conf. on Intelligent Robots and Systems*, pages 1769–1776. IEEE, 2020. 3, 4
- [17] Xingyu Jiang, Jiangwei Ren, Zizhuo Li, Xin Zhou, Dingkan Liang, and Xiang Bai. Minima: Modality invariant image matching. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2025. 2
- [18] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proc. of IEEE Intl. Conf. on Computer Vision*, pages 4015–4026, 2023. 3
- [19] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proc. of Intl. Conf. on Machine Learning*, pages 19730–19742. PMLR, 2023. 4
- [20] Shijie Lian, Hua Li, Runmin Cong, Suqi Li, Wei Zhang, and Sam Kwong. Watermask: Instance segmentation for underwater imagery. In *Proc. of IEEE Intl. Conf. on Computer Vision*, pages 1305–1315, 2023. 3, 4, 6
- [21] Shijie Lian, Ziyi Zhang, Hua Li, Wenjie Li, Laurence Tianruo Yang, Sam Kwong, and Runmin Cong. Diving into underwater: Segment anything model guided underwater salient instance segmentation and a large-scale dataset. In *Proc. of Intl. Conf. on Machine Learning*, 2024. 2, 3, 4, 6
- [22] Zhengyao Lv, Yuxiang Wei, Wangmeng Zuo, and Kwan-Yee K Wong. Place: Adaptive layout-semantic fusion for semantic image synthesis. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 9264–9274, 2024. 3, 4
- [23] Quang Nguyen, Truong Vu, Anh Tran, and Khoi Nguyen. Dataset diffusion: Diffusion-based synthetic data generation for pixel-level semantic segmentation. In *Proc. of Advances in Neural Information Processing Systems*, 2024. 2, 3
- [24] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *Proc. of Intl. Conf. on Machine Learning*, pages 16784–16804. PMLR, 2022. 3
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proc. of Intl. Conf. on Machine Learning*, pages 8748–8763. PMLR, 2021. 4

- [26] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. [4](#)
- [27] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. [3](#)
- [28] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. In *Proc. of Intl. Conf. on Learning Representations*, 2025. [3](#)
- [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. [2](#), [3](#), [4](#)
- [30] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *Proc. of Advances in Neural Information Processing Systems*, 2017. [3](#), [7](#)
- [31] Yibo Wang, Ruiyuan Gao, Kai Chen, Kaiqiang Zhou, Yingjie Cai, Lanqing Hong, Zhenguo Li, Lihui Jiang, Dit-Yan Yeung, Qiang Xu, et al. Detdiffusion: Synergizing generative and perceptive models for enhanced data generation and perception. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 7246–7255, 2024. [2](#), [3](#)
- [32] Weijia Wu, Yuzhong Zhao, Hao Chen, Yuchao Gu, Rui Zhao, Yefei He, Hong Zhou, Mike Zheng Shou, and Chunhua Shen. Datasetdm: Synthesizing data with perception annotations using diffusion models. *Proc. of Advances in Neural Information Processing Systems*, 36:54683–54695, 2023. [3](#)
- [33] Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou, and Chunhua Shen. Diffumask: Synthesizing images with pixel-level annotations for semantic segmentation using diffusion models. In *Proc. of IEEE Intl. Conf. on Computer Vision*, pages 1206–1217, 2023. [3](#)
- [34] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Proc. of Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. [7](#)
- [35] 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. [6](#), [7](#), [11](#)
- [36] Han Xue, Zhiwu Huang, Qianru Sun, Li Song, and Wenjun Zhang. Freestyle layout-to-image synthesis. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 14256–14266, 2023. [3](#), [4](#)
- [37] Tianyu Yan, Zifu Wan, Xinhao Deng, Pingping Zhang, Yang Liu, and Huchuan Lu. Mas-sam: Segment any marine animal with aggregated features. In *Proc. of Intl. Joint Conf. on Artificial Intelligence*, 2024. [3](#)
- [38] Lihe Yang, Xiaogang Xu, Bingyi Kang, Yinghuan Shi, and Hengshuang Zhao. Freemask: Synthetic images with dense annotations make stronger segmentation models. *Proc. of Advances in Neural Information Processing Systems*, 36, 2023. [2](#), [3](#), [7](#)
- [39] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. In *Proc. of Advances in Neural Information Processing Systems*, 2024. [4](#)
- [40] Hanrong Ye, Jason Kuen, Qing Liu, Zhe Lin, Brian Price, and Dan Xu. Seggen: Supercharging segmentation models with text2mask and mask2img synthesis. *arXiv preprint arXiv:2311.03355*, 2023. [3](#)
- [41] 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. [2](#), [6](#), [11](#)
- [42] 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. [2](#), [3](#), [6](#), [7](#), [11](#)
- [43] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proc. of IEEE Intl. Conf. on Computer Vision*, pages 3836–3847, 2023. [2](#), [3](#)
- [44] Pingping Zhang, Tianyu Yan, Yang Liu, and Huchuan Lu. Fantastic animals and where to find them: Segment any marine animal with dual sam. In *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, pages 2578–2587, 2024. [3](#)

A. Visualization

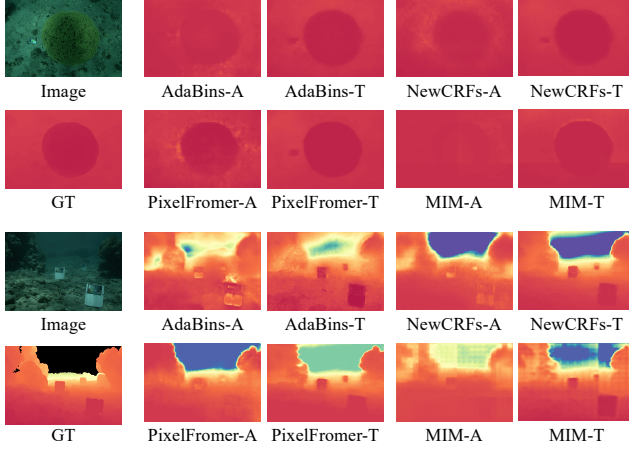


Figure A1. Qualitative results on the Sea-thru dataset [3]. ‘-A’ and ‘-T’ denote models trained on Atlantis [42] and Our SynTIDE dataset, respectively. The depth estimation results are notably improved after training on our dataset. Due to the original ‘Image’ being extremely dim, the content is hardly 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.

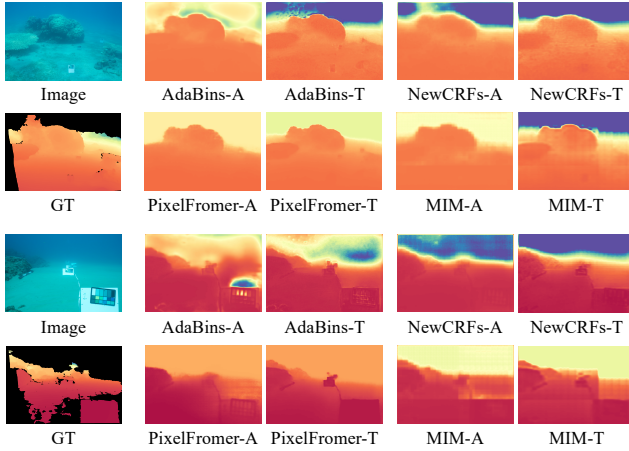


Figure A2. Qualitative results on SQUID dataset [4]. ‘-A’ and ‘-T’ denote models trained on Atlantis [42] and Our SynTIDE dataset, respectively. The depth estimation results are notably improved after training on our dataset.

A.1. Qualitative results

Fig. A1 and Fig. A2 showcase qualitative comparisons with Atlantis on the D3 and D5 subsets of the Sea-thru [3] dataset and the SQUID [4] dataset. All models trained on the SynTIDE dataset, including AdaBins [5], NeWCRFs [41], PixelFormer [1], and MIM [35], 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.

A.2. Zero-shot underwater depth data generation

Thanks to our training strategy, which fine-tunes the pre-trained text-to-image model [6] using LoRA [15] 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 capability is particularly challenging for Atlantis [42], which requires the depth map in advance as a condition.

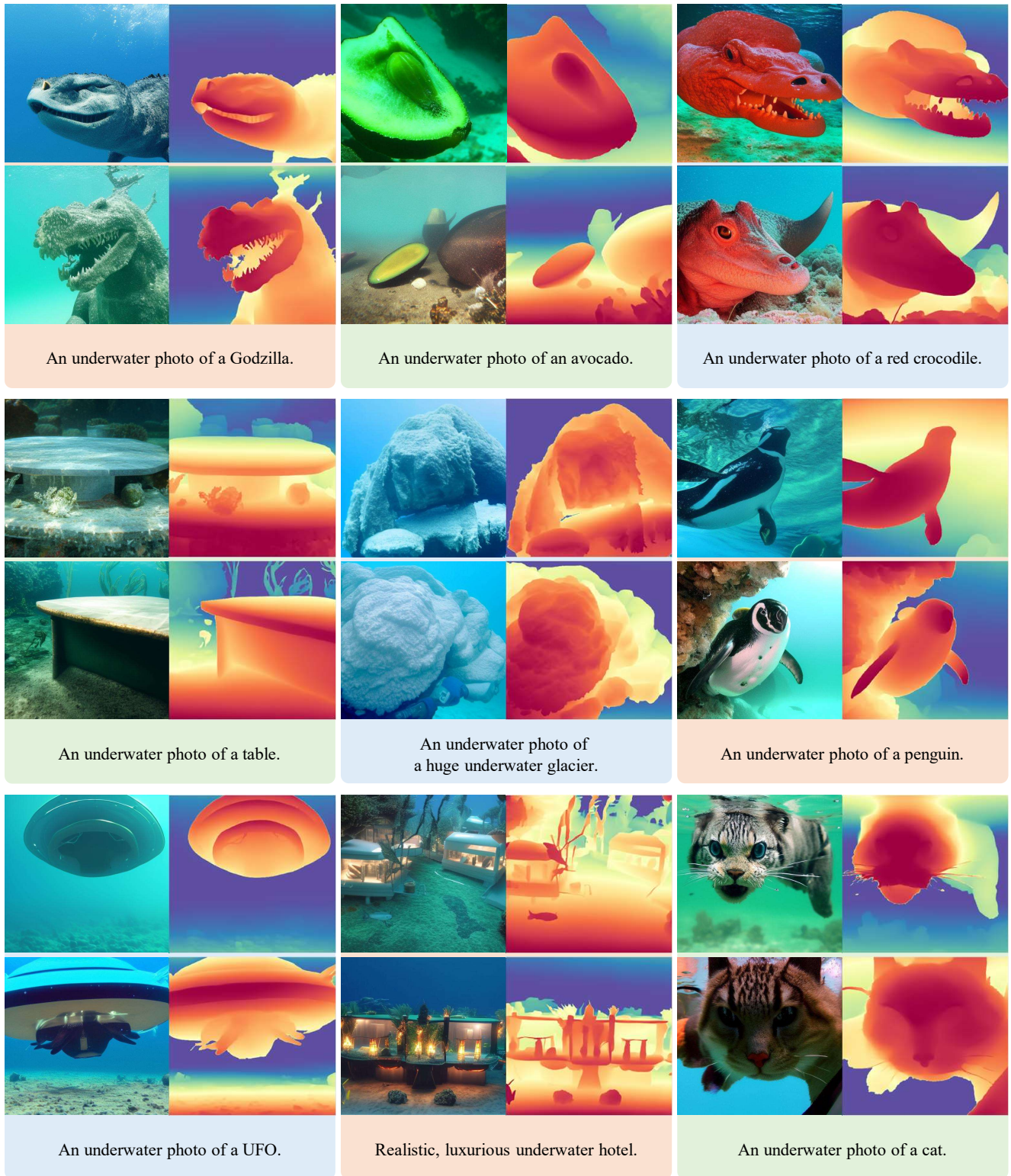


Figure A3. Representative zero-shot image-depth pairs synthesized by TIDE present strong consistency, diversity, and generalization. Images of relevant categories are not included in the training data.