

AquaticCLIP: A Vision-Language Foundation Model for Underwater Scene Analysis

Basit Alawode, Iyyakutti Iyappan Ganapathi, Sajid Javed, Naoufel Werghi, *IEEE Senior Member*, Mohammed Bennamoun, and Arif Mahmood

Abstract—The preservation of aquatic biodiversity is critical in mitigating the effects of climate change. Aquatic scene understanding plays a pivotal role in aiding marine scientists in their decision-making processes. In this paper, we introduce AquaticCLIP, a novel contrastive language-image pre-training model tailored for aquatic scene understanding. AquaticCLIP presents a new unsupervised learning framework that aligns images and texts in aquatic environments, enabling tasks such as segmentation, classification, detection, and object counting. By leveraging our large-scale underwater image-text paired dataset without the need for ground-truth annotations, our model enriches existing vision-language models in the aquatic domain. For this purpose, we construct a 2 million underwater image-text paired dataset using heterogeneous resources including YouTube, Netflix, NatGeo, etc. To fine-tune AquaticCLIP, we propose a prompt-guided vision encoder that progressively aggregates patch features via learnable prompts, while a vision-guided mechanism enhances the language encoder by incorporating visual context. The model is optimized through a contrastive pre-training loss to align visual and textual modalities. AquaticCLIP achieves notable performance improvements in zero-shot settings across multiple underwater computer vision tasks, outperforming existing methods in both robustness and interpretability. Our model sets a new benchmark for vision-language applications in underwater environments. The code and dataset for AquaticCLIP are publicly available on GitHub at [xxx](#).

Index Terms—Vision language model, Underwater scene analysis, underwater object detection, object segmentation, and object counting.

I. INTRODUCTION

Global aquatic¹ ecosystems are under severe threats from human activities such as overfishing and coastal development, along with climate change impacts [37], [49], [60], [144]. Effective conservation efforts depend on precise monitoring, which requires an accurate and automatic aquatic scene understanding system [39], [46]. However, the complexity of understanding aquatic environments demands significant expertise from ocean scientists and marine biologists, creating challenges for efficient monitoring [102], [140].

Recently, Vision-Language Models (VLMs) have gained increasing attention in the computer vision field [85], [90], [97], [103], [126]. These models are typically pre-trained using large-scale image-text paired data, readily available online

B. Alawode, I.I. Ganapathi, S. Javed, and Naoufel Werghi are with the department of computer science, Khalifa University of Science and Technology, P.O Box: 127788, Abu Dhabi, UAE. (email: sajid.javed@ku.ac.ae).

A. Mahmood is with the Department of Computer Science, Information Technology University, Lahore, Pakistan.

¹Aquatic is a broad term encompassing all water-based environments, including both freshwater and saltwater habitats. It encompasses a vast range of ecosystems, from rivers and lakes to oceans and coral reefs.

[97], [109], [134]. Pre-trained VLMs have been successfully applied to downstream tasks such as image classification [91], detection [146], tracking [72], and human action recognition [112]. Their success is primarily attributed to contrastive pre-training loss, which pulls paired images and texts closer while pushing unrelated ones apart in the embedding space [73], [97]. A notable example is Contrastive Language-Image Pre-training (CLIP), which captures rich vision-language correspondence and enables zero-shot predictions by matching the embeddings of images and texts [97], [134].

While significant progress has been made in extending CLIP to various computer vision tasks, few works have applied VLMs to aquatic environments [139], [140], [149]. For example, MarineGPT, pre-trained on 5M images, was introduced for marine image-question answering tasks [140]. More recently, MarineInst, a marine foundational model, has been proposed for segmentation and caption generation tasks [149]. However, the development of marine VLMs has lagged behind terrestrial VLMs due to the unique challenges of aquatic environments. Existing open-air datasets like COCO [78] cannot be used to train aquatic VLMs, and large-scale, domain-specific aquatic image-text pairs are not readily available. This scarcity of data makes pre-training aquatic VLMs particularly challenging.

To bridge this gap, we propose a large-scale aquatic image-text paired dataset comprising 2 million image-text pairs. We then introduce AquaticCLIP, a model that efficiently aligns aquatic images and texts for several downstream tasks, as shown in Fig. 1. Our dataset is collected from publicly available resources such as National Geographic (NatGeo) [32], [36], [40], [87], [89], [113], [114], aquatic biology textbooks and journals, YouTube aquatic documentary videos, Fishes of Australia [86], [104], [106], [118], Marine Twitter, Netflix, and the Corals of the World [119]. To our knowledge, no such paired dataset exists for aquatic scenes, except MarineInst [149], which contains only images. To further enrich the textual descriptions, we generate additional descriptions for aquatic images using MarineGPT [140] both at the image level and the instance level. For instance-level descriptions, we detect objects within an aquatic image using a pre-trained object detector and then MarineGPT is employed for each instance. These additional descriptions are then refined using a custom cleaning module to remove irrelevant keywords while retaining those relevant to the aquatic imagery.

For our AquaticCLIP model, we introduced two lightweight learnable encoders for the image and text branches, respectively. The key idea behind these encoders is to effectively transfer the VLM into the aquatic domain by leveraging the

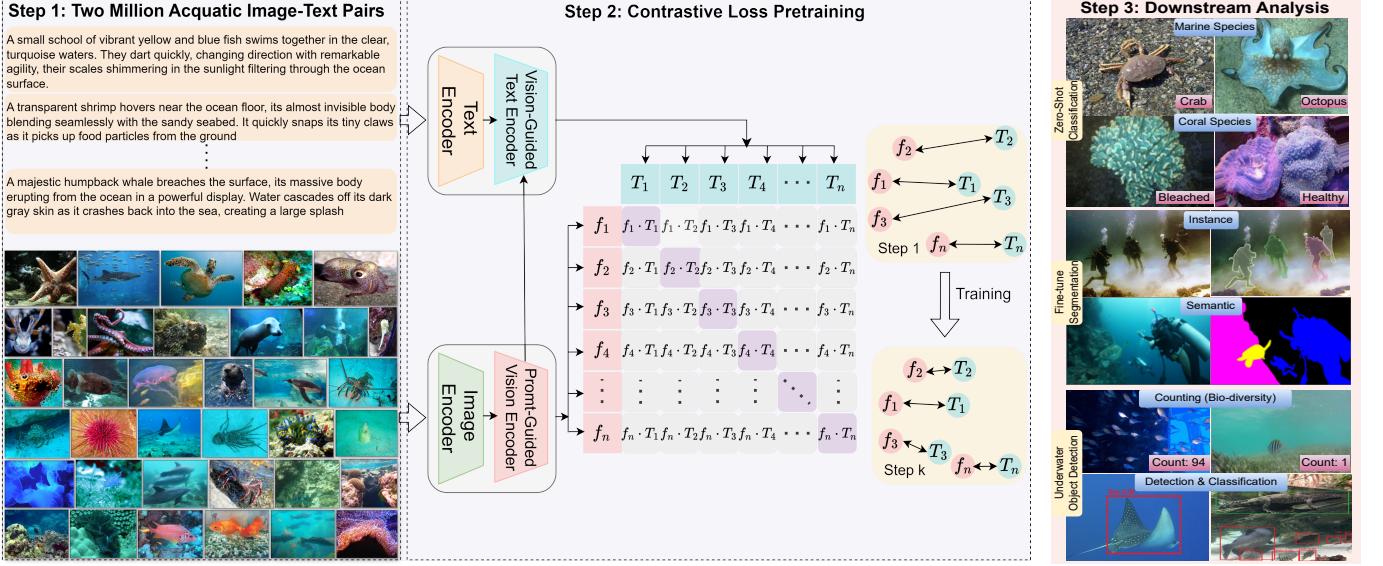


Fig. 1: (a) **Step 1: Two Million Aquatic Image-Text Pairs.** The dataset consists of paired aquatic images and enriched textual descriptions, which serve as input to the model. (b) **Step 2: Contrastive Loss Pretraining.** Text and image pairs are processed by a text encoder and an image encoder. The embeddings are aligned through contrastive loss, reducing the distance between matching pairs and improving the model’s ability to associate images with their corresponding textual descriptions. (c) **Step 3: Downstream Analysis.** AquaticCLIP performance is evaluated across various tasks such as zero-shot marine species classification, fine-tuned instance and semantic segmentation, object detection, and biodiversity counting in underwater imagery.

prior knowledge of the existing MarineGPT. To enable the VLM to process aquatic images more efficiently, we designed a prompt-guided vision encoder based on prompt-based learning. Specifically, for the image branch, we aggregate all patch features using learnable weights. A set of learnable prompt vectors is introduced to progressively guide the fusion of patch features, grouping similar ones together. This method allows each prompt to capture more global contextual information for final similarity computation.

For the text branch, we propose a vision-guided text encoder that integrates corresponding image information into the text encoder. By employing a multi-modal text encoder for guidance, knowledge is more effectively transferred to the VLM’s text encoder. The visual representations and textual descriptions learned by these two encoders are then aligned using a contrastive pre-training loss, similar to CLIP [97]. The main goal is to bring similar visual and textual concepts of aquatic images closer and push dissimilar ones apart.

We conducted extensive experimental evaluations of the AquaticCLIP model on various downstream tasks, including zero-shot and fine-tuning scenarios. Zero-shot tasks included aquatic species recognition, fine-grained fish classification, coral species classification, and cross-modal retrieval. Fine-tuning tasks involved coral segmentation, instance segmentation of aquatic imagery, semantic segmentation, aquatic object detection and classification, and aquatic object counting. Our results demonstrate significant performance improvements compared to existing State-Of-The-Art (SOTA) methods in aquatic settings.

The key contributions of this work include:

- 1) We propose a large-scale dataset of 2 million image-text pairs for aquatic VLM pre-training, with enriched textual descriptions generated by MarineGPT in zero-shot

settings at both the image and instance levels and refined to exclude undesired keywords, enhancing data quality (Sec. III-A-III-C).

- 2) We introduce a dual-encoder approach with a prompts-guided image encoder to suppress irrelevant image regions and a vision-guided text encoder for improved alignment of visual and textual representations (Sec. III-D-III-F).
- 3) Extensive evaluations on diverse aquatic vision tasks show significant gains over existing SOTA methods (Sec. IV).

The remainder of this paper is organized as follows: Sec. II reviews related work. Sec. III presents our proposed dataset and the AquaticCLIP. Sec. IV presents the experiments and results, and Sec. V concludes our work.

II. RELATED WORK

Numerous deep learning-based approaches have been proposed for various aquatic scene understanding tasks [45], [74], [102], [121], [133] including underwater image enhancement [34], [57], [123], species classification [70], [125], underwater object recognition [59], [64], coral segmentation [139], Object counting [111], and underwater tracking [133].

1. CNN-based Methods: For underwater image enhancement and restoration, many CNN-based approaches have been proposed including WaterGAN [69], Image-2-Image Translation [29], AquaGAN [34], CycleGAN [84], AGCycleGAN [123] and others [54]. For underwater object detection and classification, several methods and datasets have led to significant progress [59], [64], [130]. Khan *et al.* introduced FishNet, a benchmark dataset for fish recognition and functional trait prediction [64].

Species classification is crucial for underwater monitoring and conservation and a number of methods have been proposed in this directions [94], [96], [108]. Coral reefs, known as the “rainforests of the sea”, play a critical role in biodiversity. However, coral segmentation is a challenging task due to complex underwater visual conditions. Many methods are recently proposed for this task [141], [148]. Object counting is essential for marine biology, fisheries management, and environmental monitoring. Most approaches focus on fish and coral counting such as [18], [20], [88]. Underwater visual tracking has also seen significant advances using CNN-based trackers and new datasets [21], [50], [68], [93].

2. Transformer-based Methods: Many ViT-based techniques for detection, segmentation, counting, tracking, and classification have been proposed for underwater scenes [15], [95], [111], [127], [133]. For instance, UShape transformer [95] and transformer-driven GAN [117] have been proposed for underwater image restoration tasks. Alawode *et al.* proposed a combined approach for underwater image enhancement and visual tracking [16], and Zhang *et al.* developed a large-scale benchmark for advancing underwater object tracking [133]. A dense object counter [111] and density-guided attention [127] methods are proposed for underwater object counting task. For fish detection and classification [80], [81], and coral reef classification [15], [105] tasks, ViT-based models have also recently been proposed. A token-based selective ViT [107] and cascaded attention [136] models are proposed for fine-grained marine species classification. Additionally, some self-supervised learning methods have also been utilized for underwater image analysis [55], [101].

3. Vision-Language Models (VLMs): In the context of safeguarding aquatic biodiversity, there is an increasing need for VLMs to facilitate AI-based aquatic scene understanding systems. However, VLMs have only been sparsely applied to aquatic scene analysis [139], [140], [149]. For instance, Zheng *et al.* introduced CoralSCOP, a foundational model for the automatic segmentation of coral reefs [139], and MarineGPT, a multimodal large language model for marine object classification and question-answering [140]. Zheng *et al.* further evaluated GPT-4V for marine images, but found its performance unsatisfactory for domain-specific needs of marine biologists [138]. Recently, the MarineInst foundational model was proposed, pre-trained on 2M images for segmenting and captioning marine imagery [149]. To the best of our knowledge, VLMs have not been thoroughly explored for aquatic scene understanding, except for MarineGPT and MarineInst. Our work is the first to introduce AquaticCLIP, with comprehensive analysis and comparisons to existing SOTA methods.

III. PROPOSED AQUATICCLIP MODEL

This work introduces the Aquatic Contrastive Language-Image Pre-training (AquaticCLIP) model, which leverages a collection of 2M aquatic image-text pairs curated from several heterogeneous resources. The ground truth descriptions are also enriched by harnessing the existing VLM MarineGPT. The primary objective is to pre-train AquaticCLIP using

diverse aquatic data sourced from various platforms, enhancing its capability for zero-shot transfer across different aquatic imagery tasks. This is particularly beneficial for recognizing unfamiliar marine species and coral reef categories that were not encountered during training. Fig. 2 illustrates the key components of the proposed model, which include the construction of the 2M aquatic image-text paired dataset further enriched by unsupervised generated textual descriptions, caption cleaning, a lightweight prompt-guided vision encoder, a vision-guided text encoder, and the pre-training process. The contrastive learning approach aligns positive image-text pairs while separating negative ones. The details of these processes are discussed in the subsequent sections.

A. Aquatic Dataset Construction and Curation

Our dataset construction pipeline involves two main steps: gathering image-text paired data from multiple sources and cleaning and filtering the collected data. We assembled a dataset of 2 million aquatic image-text pairs from various resources, including YouTube videos, marine and ocean sciences textbooks and articles, the Corals of the World [119], Fishes of Australia [86], [104], [106], [118], Marine Twitter, Netflix, and National Geographic (NatGeo) [32], [36], [40], [87], [89], [113], [114].

For YouTube videos, we searched using keywords such as “underwater world”, “marine documentary”, “deep oceans”, “great barrier reef”, “aquatic scenes”, and “coral reefs” etc. For Netflix videos, we explored hundreds of documentaries, including “My Octopus Teacher”, “Last Breath”, and “Wonders of the Reef”, etc. Subtitles provided by both resources were used to generate aligned image-text pairs, which were manually checked and refined. Unique frames were extracted every 50 seconds from the videos, which often contained challenges like low visibility, motion blur, background clutter, and color distortions.

Additionally, We utilized 1200 diverse textbooks on marine biology and oceanography, along with research articles from ocean and marine journals and NatGeo magazines, to further enrich the dataset. Figures and captions were extracted using PDF-Figures 2.0 tool [30], and we manually refined the data to ensure the selected images had meaningful associated text. Images not related to aquatic environments were discarded. We also included image-text pairs from the Corals of the World repository and Fishes of Australia, only selecting pairs with detailed descriptions. Furthermore, we used the Twitter platform to search for relevant content using hashtags like #MarineBiology, #Oceans, and #Fisheries, considering only channels with over 100 followers. After a thorough cleaning and filtering process, we retained 2 million high-quality image-text pairs representing a diverse range of aquatic scenes. *More details are provided in the supplementary material.*

B. Unsupervised Generation of Image and Instance-Level Descriptions (Fig. 2 (b))

In order to enrich the ground-truth textual descriptions, we generated additional textual descriptions at both the image and instance levels using a VLM MarineGPT [140], which

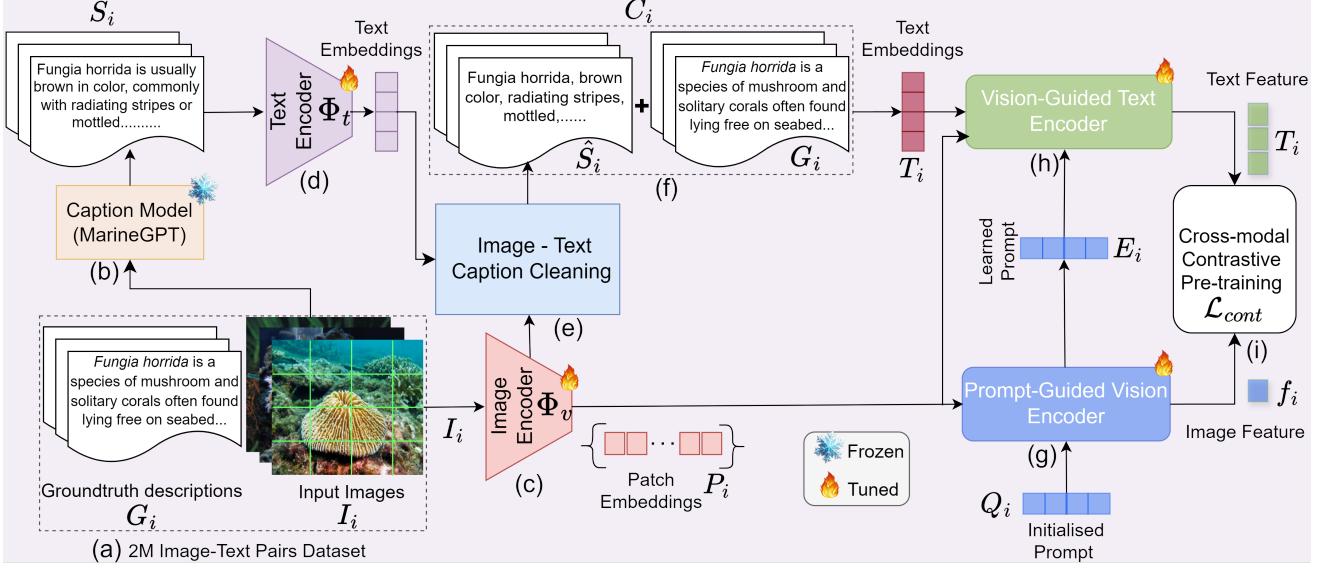


Fig. 2: **Overview of AquaticCLIP architecture and training process.** (a) Shows a set of input image-text pairs. (b) A caption model (MarineGPT) generates textual descriptions for the images. (c) Input images are divided into patches and processed by the image encoder Φ_v to produce patch embeddings P_i . (d) The generated textual descriptions S_i are processed by the text encoder Φ_t to produce text embeddings. (e)-(f) The textual description S_i is then cleaned by an image-text caption cleaning module to produce refined descriptions \hat{S}_i which are then combined with groundtruth descriptions G_i to produce enriched textual description data C_i . Both image and text embeddings are refined using (h) vision-guided text encoding and (g) prompt-guided vision encoding. The learned prompts E_i guide the fusion of patch embeddings, while initialized prompts Q_i are used to enhance the visual representation. (i) The final image and text features are aligned using a cross-modal contrastive pre-training loss \mathcal{L}_{cont} , ensuring a stronger association between text and image representations.

includes a frozen ViT image encoder and Q-former. At the image level, each image I_i is input into MarineGPT to generate its corresponding textual descriptions.

For the instance level, we pre-trained MRegionCLIP, an object detector based on RegionCLIP [142] and MarineDet [48] and applied it to our 2M imagery dataset to detect all instances in zero-shot settings. Each instance was then passed through MarineGPT to generate a textual description. Specifically, we used the following prompt template: “The image is < image >. Describe the object in this image.”, where < image > is the image token. The generated textual descriptions S_i at the image and instance levels were then cleaned using our cleaning module, which is explained in the following subsection.

C. Semantic Filtering and Cleaning of Generated Textual Descriptions (Figs. 2 (c)-(f))

The generated textual descriptions S_i may contain noise, such as broken sentences, incorrect descriptions, or irrelevant keywords. To address these issues, we developed a textual description cleaning module aimed at identifying the semantically closest and most relevant keywords.

In this process, each generated textual description S_i is broken down into a set of k -keywords ($\{s_i^j\}_{j=1}^k$). For each keyword, we compute its cosine similarity with the image embedding as follows:

$$\hat{S}_i = \operatorname{argmax}_{s \in S} \langle \Phi_v(I_i) \cdot \Phi_t(s_i^j) \rangle, \quad (1)$$

where Φ_v is a vision encoder followed by an MLP, and Φ_t is a text encoder from the CLIP model. We retain the top-p%

of keywords in \hat{S}_i , discarding the rest as noise. This cleaning process ensures that the remaining keywords are semantically aligned with the visual content, improving the quality of the textual descriptions.

For each image in our dataset, we manually verified ground truth textual descriptions G_i . During the pre-training stage, the refined keywords \hat{S}_i both at the image and instance levels were combined with ground-truth descriptions G_i to generate more enriched and comprehensive textual data C_i for further processing. *Our AquaticCLIP model is pre-trained on 2M image-text pairs, where the ground-truth descriptions of the images were combined with refined keywords.*

D. Prompt-guided Vision Encoder (Fig. 3 (a))

To generate efficient visual embeddings for each aquatic image I_i , we aggregate the patch features of the input image using learned visual prompts. First, the input image I_i is divided into n_p non-overlapping patches $\{w_i^j\}_{j=1}^{n_p}$, each of size $m \times m$. These patches are fed into the pre-trained image encoder Φ_v to produce embeddings $P_i = \{\mathbf{p}_i^j\} \in \mathbb{R}^{d_p \times n_p}$.

To effectively aggregate these patch embeddings into a final image-level embedding for similarity calculation, we designed a prompt-guided image encoder, as illustrated in Fig. 3 (a). We randomly initialize a set of learnable prompt features $Q_i = \{r_i^j\}_{j=1}^{n_r} \in \mathbb{R}^{d_p \times n_r}$, where n_r represents the number of learnable prompts. These prompts guide the progressive fusion of patch embeddings. Cross-attention is then computed using the visual embeddings as keys $\mathbf{K}_i = \mathbf{P}_i$ and values $\mathbf{V}_i = \mathbf{P}_i$,

while the prompts \mathbf{Q}_i serve as queries.

$$\mathbf{E}_i = \text{Softmax} \left(\frac{\mathbf{Q}_i \mathbf{K}_i^\top}{\sqrt{d_p}} \right) \mathbf{V}_i, \mathbf{E}_i = \text{Norm}(\mathbf{E}_i) + \mathbf{Q}_i, \quad (2)$$

The learnable prompts help prioritize patches with high semantic similarity, resulting in a more meaningful image-level representation that captures global contextual information. The final image-level features are derived using an attention-based feature fusion method, as shown below:

$$\mathbf{E}'_i = \mathbf{W}_1 \mathbf{E}_i, \mathbf{e}_i = \exp(\mathbf{W}_3^\top (\tanh(\mathbf{W}_2 \mathbf{E}'_i))), \quad (3)$$

Here, $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d_p \times d_p}$ and $\mathbf{W}_3 \in \mathbb{R}^{1 \times d_p}$ are learnable matrices, and the softmax function is used to compute attention weights $\mathbf{a}_i(j)$. The image-level representation \mathbf{f}_i is then computed as follows:

$$\mathbf{a}_i(j) = \frac{\mathbf{e}_i(j)}{\sum_{k=1}^{n_r} \mathbf{e}_i(k)}, \text{ and } \mathbf{f}_i = \mathbf{W}_4 \sum_{j=1}^{n_r} \mathbf{a}_i(j) \mathbf{E}'_i(j), \quad (4)$$

where $\mathbf{W}_4 \in \mathbb{R}^{d_p \times d_p}$ is a learnable weight matrix and $\mathbf{E}'_i(j)$ is a column vector of \mathbf{E}'_i .

E. Vision-guided Text Encoder (Fig. 3 (b))

In the text encoder branch, the enriched textual descriptions \mathbf{C}_i are fed into the CLIP text encoder to obtain textual representations \mathbf{T}_i corresponding to the descriptions of the i -th image. These representations are then passed through a vision-guided attention layer for refinement. The patch features \mathbf{P}_i and the learned prompts \mathbf{E}_i are concatenated as \mathbf{V}_i , which serves as the key \mathbf{K}_t and value \mathbf{V}_t , while the textual representations \mathbf{T}_i are used as the query, as shown in (Fig. 3 (b)). The vision-guided attention mechanism is computed as follows:

$$\mathbf{U}_i = \text{Softmax} \left(\frac{\mathbf{T}_i \mathbf{K}_{t,i}^\top}{\sqrt{d_p}} \right) \mathbf{V}_{t,i}, \mathbf{T}_i = \mathbf{T}_i + \mathbf{U}_i, \quad (5)$$

This context-guided text encoder further refines the textual features by incorporating image context and learned visual prompts. This process enhances the alignment between images and texts, improving the performance of the AquaticCLIP model.

F. Cross-Modal Contrastive Loss for Vision-Language Alignment (Fig. 2 (i))

We pre-train our prompt-guided vision encoder and vision-guided text encoder using a cross-modal contrastive loss function. This loss is formulated as a temperature-scaled vision-language pre-training loss, similar to W -way classification, where W represents the batch size of image-text pairs involved in the training process [25], [97], [115]. Given a batch of W paired normalized image and text embeddings $\{\mathbf{f}_i, \mathbf{T}_i\}$, we minimize the contrastive loss in two directions: image-to-text ($i \rightarrow t$) and text-to-image ($t \rightarrow i$) as:

$$\mathcal{L}_{i2t} = -\frac{1}{W} \sum_{i=1}^W \log \frac{\exp(\tau \mathbf{T}_i^\top \mathbf{f}_i)}{\sum_{j=1}^W \exp(\tau \mathbf{T}_i^\top \mathbf{f}_j)}, \quad (6)$$

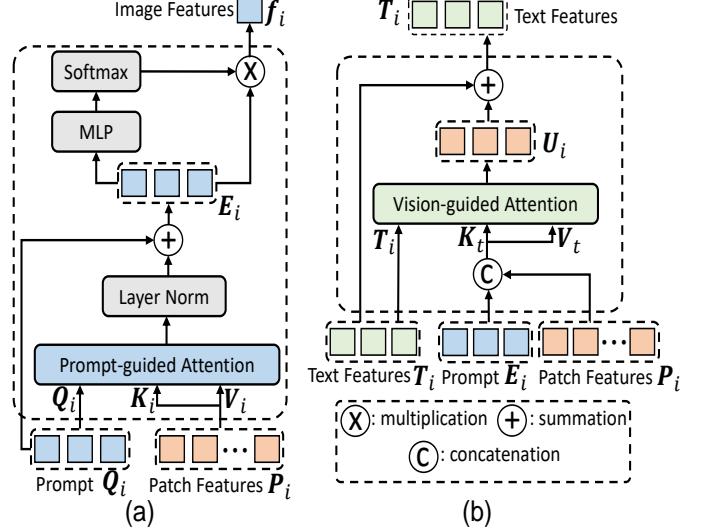


Fig. 3: (a) **Prompt-Guided Vision Encoder**: The prompt-guided attention mechanism combines patch features \mathbf{P}_i with initialized prompts \mathbf{Q}_i through layer normalization and an MLP, followed by softmax to produce the final image features \mathbf{f}_i . (b) **Vision-Guided Text Encoder**: Text embeddings \mathbf{T}_i are refined using a vision-guided attention mechanism, where patch features \mathbf{P}_i , learned prompts \mathbf{E}_i , and text embeddings \mathbf{T}_i are concatenated to compute attention \mathbf{U}_i , which further enhances \mathbf{T}_i .

$$\mathcal{L}_{t2i} = -\frac{1}{W} \sum_{j=1}^W \log \frac{\exp(\tau \mathbf{f}_j^\top \mathbf{T}_i)}{\sum_{i=1}^W \exp(\tau \mathbf{f}_j^\top \mathbf{T}_i)}, \quad (7)$$

where τ is a learnable temperature parameter that controls the smoothness of the distribution [97]. The overall contrastive loss, \mathcal{L}_{cont} , is the sum of both losses: $\mathcal{L}_{cont} = \mathcal{L}_{i2t} + \mathcal{L}_{t2i}$. This loss minimizes the distance between embeddings of positive image-text pairs and maximizes the distance between negative pairs $(\mathbf{f}_i, \mathbf{T}_j)$, where $i \neq j$, ensuring that images and texts with the same semantic content have similar representations in the feature space.

G. Zero-shot Transfer for Image Classification

Radford *et al.* introduced the concept of zero-shot classification using a prompt-based approach [97]. In our method, each class in the test dataset is converted into one or more text prompts using predefined templates, such as “An image of {Sea Urchins}.” or “An image of {Oyster}.” For each test image, we compute the ℓ_2 normalized embeddings using our prompt-guided vision encoder and vision-guided text encoder. We then calculate the cosine similarity between the test image and the set of testing prompts to find the best match, resulting in zero-shot classification. Additional details are provided in the supplementary material.

IV. EXPERIMENTAL EVALUATIONS

We conducted extensive experiments to evaluate the performance of the proposed AquaticCLIP model across various tasks, including zero-shot classification of marine species, fine-grained fish, and coral species classification. Additionally, we performed zero-shot cross-modal retrieval tasks for aquatic images (see supplementary material). For downstream tasks,

we applied fine-tuned instance segmentation, semantic segmentation of underwater imagery, as well as marine object detection, classification, and counting (*see supplementary material*). These experiments, covering a range of classification, detection, and segmentation tasks, allowed us to thoroughly assess AquaticCLIP’s performance. We also compared our results with SOTA methods, including both VLMs and vision-only approaches.

A. Training and Implementation Details

The architecture of AquaticCLIP consists of a frozen domain-specific captions generator MarineGPT [140], the CLIP [97] image encoder using the ViT-B/16-224 [38], and a transformer-based text encoder [98]. We fine-tuned four components: the image encoder, text encoder, prompt-guided vision encoder, and vision-guided text encoder, all using the cross-modal contrastive loss described in Sec. III-F. We employed the Adam optimizer [83] with an initial learning rate of 1×10^{-4} and a weight decay of 1×10^{-5} . The model was trained for 80 epochs on four A100 GPUs with a batch size of 512. We set the number of prompts to 20. For fair comparisons, we utilized the same evaluation metrics as used by the SOTA methods. For the classification tasks, we reported both accuracy and F_1 measure scores. For object detection, we reported mAP₅₀ metric. *Additional details are provided in supplementary material.*

B. Underwater Datasets and Tasks

For the zero-shot marine species classification task, we utilized two datasets: Marine Animal Images (MAI) [7] and Sea Animals Images (SAI) [8]. For zero-shot fine-grained classification, we employed three datasets: FishNet [64], FishNet Open Images (F NOI) [62], and Large-Scale Fish (LSF) [116]. For zero-shot fine-grained coral species classification, we used the Coral Species Classification (CSC) [47] and Coral Classification (CC) [35] datasets. For object detection and classification, we used four datasets including FishNet, DeepFish [100], Brackish [94], and URPC [13]. *Additional details are provided in supplementary material.*

C. Ablation Studies

We conducted ablation studies to highlight the contributions of each component of the proposed AquaticCLIP model. The evaluations were performed on zero-shot classification tasks, reporting F_1 scores across six independent external datasets (MAI, SAI, FishNet, F NOI, LSF, and CC) that were not used during the pre-training phase. **1. FrozenCLIP vs. FinetuneCLIP (Table V):** In FrozenCLIP, pre-trained image and text encoders were used to generate visual and text embeddings. In FinetuneCLIP, these encoders were fine-tuned using the contrastive pre-training loss on our 2M dataset. We observed that fine-tuning resulted in improved performance across all datasets. **2. AquaticCLIP₁ vs. AquaticCLIP (Table V):** AquaticCLIP₁ used frozen image and text encoders, with only the PGVE (Prompt-Guided Vision Encoder) and VGTE (Vision-Guided Text Encoder) fine-tuned. Significant performance improvements were observed when all four encoders were fine-tuned in AquaticCLIP. **3. Importance of**

PGVE and VGTE Components (Table V): In AquaticCLIP₂ and AquaticCLIP₄, VGTE was removed, and original textual embeddings were used, while in AquaticCLIP₃ and AquaticCLIP₅, PGVE was removed and image-level embeddings from MLP were used. Both cases showed performance reductions compared to AquaticCLIP₁ and AquaticCLIP, highlighting the importance of these components.

4. Importance of Textual Description Cleaning Module (TDCM) (Table VI):

In AquaticCLIP₆, removing the TDCM module led to decreased performance compared to the full AquaticCLIP model.

5. Importance of Instance-level Text and Image-level Text (Table VI):

In AquaticCLIP₇ and AquaticCLIP₈, either instance-level or image-level text descriptions were removed. Performance dropped in both cases compared to the proposed AquaticCLIP, which utilized both levels of text. For fine-grained classification (FishNet, F NOI, LSF), instance-level descriptions performed better, while for coarse-grained classification (MAI, SAI), image-level captions were more effective. *More ablations are provided in supplementary material.*

D. SOTA Methods for Comparison

We compared our AquaticCLIP model to a wide range of SOTA methods across various underwater image analysis tasks. For zero-shot classification, we compared the performance of our AquaticCLIP model with existing VLMs, including Frozen CLIP [97], Finetune CLIP [97], and prompt-based VLMs like CoOp [143] and MAPLE [65], as well as GPT4V [129], BLIP2 [71], and MarineGPT [140]. We finetuned CoOp and MAPLE using the original author’s provided source codes. For supervised classification tasks, we compared AquaticCLIP with models such as ResNet-34/50/101 [51], ViT-S/B/L [17], BeiT [19], ConvNeXt [82], ConvNeXt [82] + Focal Loss (FL) [77], and ConvNeXt [82] + Class-Balanced (CB) [31]. To ensure a fair comparison, we used the same settings as FishNet [64]. For object detection, we compared with FasterRCNN [99], YOLOF [24], TOOD [42], MarintInst [149], MarineDet [48]. *Further details are given in the supplementary material.*

E. Zero-shot Comparisons

Table III presents the zero-shot classification results and comparisons with SOTA VLM-based methods across seven datasets. AquaticCLIP consistently outperformed existing SOTA VLMs by a significant margin on all datasets. Significantly, on the CSC dataset, AquaticCLIP achieved a zero-shot performance of 96.80% accuracy and 96.40% F_1 score, the highest across all datasets. In more challenging fine-grained fish classification tasks on the FishNet, F NOI, and LSF datasets, AquaticCLIP achieved F_1 scores of 84.20%, 80.10%, and 93.40%, respectively, due to the incorporation of instance-level captions in the model.

F. Linear Probe Evaluations

In this experiment, we conducted a linear probe evaluation of the AquaticCLIP model and compared it against several

TABLE I: **Architectural ablation** highlighting the importance of the Prompt-Guided Vision Encoder (PGVE) and Vision-Guided Text Encoder (VGTE), as well as the impact of frozen versus fine-tuned image Φ_v and text encoders Φ_t . The results, reported as F_1 scores, reflect zero-shot classification performance on six datasets. The fully fine-tuned AquaticCLIP model with both PGVE and VGTE achieves the highest scores, demonstrating the effectiveness of fine-tuning and these key components.

Variants	Frozen (Φ_v)	Frozen (Φ_t)	PGVE	VGTE	MAI	SAI	FishNet	F NOI	LSF	CC
Frozen CLIP	✓	✓	✗	✗	0.692	0.702	0.651	0.622	0.772	0.752
AquaticCLIP ₁	✓	✓	✓	✓	0.736	0.754	0.762	0.672	0.831	0.823
AquaticCLIP ₂	✓	✓	✓	✗	0.721	0.725	0.667	0.651	0.780	0.783
AquaticCLIP ₃	✓	✓	✗	✓	0.718	0.716	0.731	0.696	0.807	0.811
Finetune CLIP	Finetune	Finetune	✗	✗	0.772	0.847	0.771	0.750	0.853	0.831
AquaticCLIP ₄	Finetune	Finetune	✓	✗	0.838	0.890	0.821	0.753	0.892	0.910
AquaticCLIP ₅	Finetune	Finetune	✗	✓	0.842	0.876	0.808	0.766	0.912	0.932
AquaticCLIP	Finetune	Finetune	✓	✓	0.871	0.923	0.842	0.801	0.934	0.953

TABLE II: **Ablation study** showing the effect of the Textual Descriptions Cleaning Module (TDCM) and the use of instance-level and image-level textual descriptions on AquaticCLIP zero-shot classification performance (F_1 scores). The full AquaticCLIP model, with TDCM and both text levels, delivers the best results across all datasets, especially on MAI (0.871), SAI (0.923), and CC (0.953). Variants with components removed or altered exhibit reduced performance, highlighting the importance of each component.

Variants	TDCM	Instance Text	Image Text	MAI	SAI	FishNet	F NOI	LSF	CC
AquaticCLIP	✓	✓	✓	0.871	0.923	0.842	0.801	0.934	0.953
AquaticCLIP ₆	✗	✓	✓	0.854	0.891	0.804	0.786	0.897	0.934
AquaticCLIP ₇	✓	✗	✓	0.853	0.906	0.804	0.765	0.911	0.933
AquaticCLIP ₈	✓	✓	✗	0.840	0.892	0.823	0.784	0.921	0.915

TABLE III: Zero-shot and supervised classification performance comparison in terms of accuracy and F_1 score of AquaticCLIP against SOTA VLMs and vision-only models. AquaticCLIP consistently outperforms all models across multiple datasets, excelling in both zero-shot and supervised tasks. It achieves top F_1 scores, particularly in MAI, SAI, FishNet, and CC datasets, demonstrating superior generalization and classification accuracy compared to traditional vision-only models and other VLMs. Family classification performance is reported for FishNet, while CSC is excluded from supervised methods due to its small size.

Zero-Shot VLMs	MAI [7]	SAI [8]	FishNet [64]	F NOI [62]	LSF [16]	CSC [47]	CC [35]
Frozen CLIP [97]	0.702 0.692	0.711 0.702	0.663 0.651	0.642 0.622	0.770 0.772	0.809 0.783	0.763 0.752
Finetune CLIP	0.802 0.772	0.856 0.847	0.770 0.771	0.763 0.750	0.864 0.853	0.862 0.841	0.845 0.831
CoOp [143]	0.853 0.822	0.866 0.853	0.752 0.744	0.764 0.752	0.863 0.854	0.903 0.888	0.868 0.866
MAPLE [65]	0.861 0.834	0.867 0.860	0.750 0.748	0.774 0.769	0.864 0.859	0.903 0.893	0.881 0.876
GPT4V [129]	0.831 0.811	0.832 0.834	0.801 0.791	0.758 0.743	0.892 0.881	0.854 0.841	0.881 0.876
BLIP2 [71]	0.813 0.801	0.821 0.818	0.783 0.788	0.728 0.727	0.893 0.880	0.793 0.782	0.862 0.853
MarineGPT [140]	0.862 0.844	0.892 0.883	0.823 0.815	0.776 0.769	0.912 0.905	0.918 0.903	0.881 0.876
AquaticCLIP ₁	0.766 0.736	0.746 0.754	0.772 0.762	0.684 0.672	0.844 0.831	0.897 0.882	0.835 0.823
AquaticCLIP ₇	0.879 0.853	0.912 0.906	0.814 0.804	0.773 0.765	0.917 0.911	0.944 0.932	0.938 0.933
AquaticCLIP ₈	0.855 0.840	0.898 0.892	0.821 0.823	0.796 0.784	0.932 0.921	0.942 0.938	0.926 0.915
AquaticCLIP	0.892 0.871	0.935 0.923	0.850 0.842	0.822 0.801	0.942 0.934	0.968 0.964	0.961 0.953
Supervised Vision Models	MAI [7]	SAI [8]	FishNet [64]	F NOI [62]	LSF [16]	CSC [47]	CC [35]
ResNet-34 [51]	0.821 0.802	0.731 0.726	0.408 0.423	0.475 0.456	0.761 0.756	-	0.833 0.821
ResNet-50 [51]	0.830 0.811	0.739 0.728	0.403 0.428	0.504 0.488	0.773 0.764	-	0.847 0.840
ResNet-101 [51]	0.842 0.816	0.745 0.740	0.363 0.354	0.511 0.496	0.809 0.792	-	0.872 0.861
ViT-S [17]	0.862 0.856	0.783 0.774	0.379 0.367	0.652 0.633	0.844 0.833	-	0.891 0.883
ViT-B [17]	0.880 0.873	0.827 0.812	0.429 0.412	0.671 0.665	0.881 0.876	-	0.910 0.903
ViT-L [17]	0.893 0.881	0.856 0.833	0.484 0.476	0.714 0.706	0.914 0.902	-	0.919 0.920
BeiT [19]	0.881 0.873	0.856 0.845	0.542 0.522	0.704 0.688	0.871 0.867	-	0.895 0.883
ConvNeXt [82]	0.847 0.852	0.812 0.803	0.606 0.587	0.714 0.702	0.827 0.807	-	0.914 0.902
ConvNeXt [82]+ FL [77]	0.881 0.870	0.834 0.822	0.551 0.544	0.738 0.722	0.831 0.842	-	0.905 0.903
ConvNeXt [82]+ CB [31]	0.883 0.872	0.851 0.842	0.613 0.605	0.745 0.731	0.870 0.855	-	0.917 0.905
AquaticVision (Linear Probing)	0.912 0.890	0.945 0.924	0.922 0.902	0.846 0.833	0.942 0.934	-	0.947 0.931
AquaticCLIP (Linear Probing)	0.915 0.893	0.951 0.944	0.934 0.923	0.882 0.867	0.968 0.961	-	0.963 0.958

SOTA supervised vision-only models. For this purpose, the vision encoder of the AquaticCLIP model is kept frozen and only a linear classifier is trained on top of the pre-extracted visual representations. We also pre-trained our vision-only model, AquaticVision, using contrastive loss on our 2M images dataset in a self-supervised learning setting. For this purpose, DINOv2 pre-training paradigm is utilized based on ViT architecture [92].

The supervised comparisons between AquaticCLIP and exist-

ing SOTA vision-only models are shown in Table III. AquaticCLIP outperformed SOTA vision-based methods, achieving F_1 scores of 92.30%, 86.70%, and 96.10% on fine-grained classification datasets. We observed that AquaticCLIP also outperforms the AquaticVision model. While AquaticVision ranked as the second-best performer, it significantly outperformed other vision-only models due to its contrastive pre-training in a self-supervised learning manner on the 2M aquatic images.

TABLE IV: Object detection and classification results (mAP₅₀).

Methods	FishNet	DeepFish	Brackish	URPC
FasterRCNN [99]	0.284	0.814	0.788	0.475
YOLOF [24]	0.672	0.806	0.813	0.511
TOOD [42]	0.811	0.766	0.805	0.507
MRegionCLIP	<u>0.867</u>	0.855	<u>0.842</u>	0.758
MarineInst [149]	0.868	0.854	0.841	0.779
MarineDet [48]	-	-	-	0.706
AquaticDet	0.903	0.891	0.877	0.837

G. Object Detection and Classification Results

In this experiment, we replaced ResNet-50 in MRegionCLIP with our pre-trained image encoder Φ_v and applied the same fine-tuning settings as discussed in the supplementary. We named this model AquaticDet and compared it against SOTA methods. Table IV shows object detection and classification results across four different datasets, with AquaticDet achieving the best mAP₅₀ scores across all compared datasets by a significant margin. This superior performance is attributed to pre-training on the 2M image-text paired dataset, which allowed AquaticDet to extract highly efficient and effective visual features. The inclusion of the prompt-guided vision encoder and vision-guided text encoder, along with comprehensive captions at both the image and instance levels, contributed to substantial performance improvements, even under challenging conditions.

H. Computational Complexity

We also evaluated the computational complexity of AquaticCLIP during the inference stage across various tasks, including zero-shot classification, supervised object detection and classification, supervised segmentation, and object counting. For zero-shot classification, AquaticCLIP took 0.80 seconds, AquaticSAM took 1.23 seconds, AquaticDet took 1.19 seconds, while the AquaticOC model took 1.31 seconds to count objects. We used the same hardware settings as discussed above.

V. CONCLUSION

Current aquatic and underwater VLMs rely on paired image-text data for pre-training. In this work, we introduced AquaticCLIP, pre-trained using real aquatic image-text pairs and additional generated textual descriptions at both image and instance levels. To achieve this, we built a 2M image-text paired dataset sourced from various online repositories. We proposed a novel vision-language alignment model where the vision encoder is guided by learned prompts, and the text encoder benefits from visual prompts. Both lightweight encoders are pre-trained using cross-modal contrastive supervision for enhanced vision-language alignment. AquaticCLIP was evaluated across a diverse set of marine vision tasks, including zero-shot fine-grained object classification, fine-tuned instance and semantic segmentation, and object detection, and counting. Our model consistently delivered superior results compared to existing SOTA methods designed specifically for marine environments, demonstrating its robustness and effectiveness across multiple aquatic vision tasks.

REFERENCES

- [1] “,” Available on Reelfile survey.
- [2] “,” Available on Reeflex.
- [3] “Aquarium dataset,” Available on Aquarium.
- [4] “Flicker Images.”
- [5] “Getty images.”
- [6] “Hk reef fish Images.”
- [7] “Marine Animal Images,” Available on Kaggle.
- [8] “Sea Animals Image Dataset,” Available on Kaggle.
- [9] “Shutterstock Image.”
- [10] “Underwater trash detection dataset,” Available on Roboflow.
- [11] Available on EOL, 2018.
- [12] “Ozfish dataset - machine learning dataset for baited remote underwater video stations,” 2020.
- [13] “URPC dataset,” Available on Kaggle, 2020.
- [14] “Oceanic life dataset,” Available on Kaggle, 2023.
- [15] B. Ai, X. Liu, Z. Wen, L. Wang, H. Ma, and G. Lv, “A novel coral reef classification method combining radiative transfer model with deep learning,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2024.
- [16] B. Alawode, Y. Guo, M. Ummar, N. Werghi, J. Dias, A. Mian, and S. Javed, “Utb180: A high-quality benchmark for underwater tracking,” in *Proceedings of the Asian Conference on Computer Vision*, 2022, pp. 3326–3342.
- [17] D. Alexey, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv: 2010.11929*, 2020.
- [18] K. M. Babu, D. Bentall, D. T. Ashton, M. Puklowski, W. Fantham, H. T. Lin, N. P. Tuckey, M. Wellenreuther, and L. K. Jesson, “Computer vision in aquaculture: a case study of juvenile fish counting,” *Journal of the Royal Society of New Zealand*, vol. 53, no. 1, pp. 52–68, 2023.
- [19] H. Bao, L. Dong, S. Piao, and F. Wei, “Beit: Bert pre-training of image transformers,” *arXiv preprint arXiv:2106.08254*, 2021.
- [20] A. B. Burguera, F. Bonin-Font, D. Chatzievangelou, M. V. Fernandez, and J. Aguzzi, “Deep learning for detection and counting of nephrops norvegicus from underwater videos,” *ICES Journal of Marine Science*, p. fsae089, 2024.
- [21] L. Cai, N. E. McGuire, R. Hanlon, T. A. Mooney, and Y. Girdhar, “Semi-supervised visual tracking of marine animals using autonomous underwater vehicles,” *International Journal of Computer Vision*, vol. 131, no. 6, pp. 1406–1427, 2023.
- [22] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, “Emerging properties in self-supervised vision transformers,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 9650–9660.
- [23] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 834–848, 2018.
- [24] Q. Chen, Y. Wang, T. Yang, X. Zhang, J. Cheng, and J. Sun, “You only look one-level feature,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 13 039–13 048.
- [25] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [26] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, “Masked-attention mask transformer for universal image segmentation,” in *CVPR*, 2022.
- [27] Y. Cheng, J. Zhu, M. Jiang, J. Fu, C. Pang, P. Wang, K. Sankaran, O. Onabola, Y. Liu, D. Liu, and Y. Bengio, “Flow: A dataset and benchmark for floating waste detection in inland waters,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 10 953–10 962.
- [28] C. Chin, “Marine fouling images,” 2019.
- [29] Y. Cho, H. Jang, R. Malav, G. Pandey, and A. Kim, “Underwater image dehazing via unpaired image-to-image translation,” *International Journal of Control, Automation and Systems*, vol. 18, pp. 605–614, 2020.
- [30] C. Clark and S. Divvala, “Pdffigures 2.0: Mining figures from research papers,” in *Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries*, 2016, pp. 143–152.
- [31] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, “Class-balanced loss based on effective number of samples,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 9268–9277.

- [32] K. Dell, "Ocean acidification's impact on marine ecosystems," *National Geographic*, May 2024, accessed: November 9, 2024.
- [33] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPR Workshops)*. Los Alamitos, CA, USA: IEEE Computer Society, jun 2009, pp. 248–255.
- [34] C. Desai, B. S. S. Reddy, R. A. Tabib, U. Patil, and U. Mudenagudi, "Aquagan: Restoration of underwater images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 296–304.
- [35] A. Dighe, "Corals classification," Available on Kaggle.
- [36] J. Doe, "The wonders of coral reefs," *National Geographic*, vol. 245, no. 4, pp. 56–71, April 2023.
- [37] S. C. Doney, M. Ruckelshaus, J. Emmett Duffy, J. P. Barry, F. Chan, C. A. English, H. M. Galindo, J. M. Grebmeier, A. B. Hollowed, N. Knowlton *et al.*, "Climate change impacts on marine ecosystems," *Annual review of marine science*, vol. 4, no. 1, pp. 11–37, 2012.
- [38] A. Dosovitskiy, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [39] C. M. Duarte, S. Agusti, E. Barbier, G. L. Britten, J. C. Castilla, J.-P. Gattuso, R. W. Fulweiler, T. P. Hughes, N. Knowlton, C. E. Lovelock *et al.*, "Rebuilding marine life," *Nature*, vol. 580, no. 7801, pp. 39–51, 2020.
- [40] J. Evers, "The decline of australia's great barrier reef," *National Geographic*, September 2023, accessed: November 9, 2024.
- [41] Y. Fang, S. Yang, X. Wang, Y. Li, C. Fang, Y. Shan, B. Feng, and W. Liu, "Instances as queries," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 6910–6919.
- [42] C. Feng, Y. Zhong, Y. Gao, M. R. Scott, and W. Huang, "Tood: Task-aligned one-stage object detection," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE Computer Society, 2021, pp. 3490–3499.
- [43] M. Fulton, J. Hong, M. J. Islam, and J. Sattar, "Robotic detection of marine litter using deep visual detection models," in *2019 International Conference on Robotics and Automation (ICRA)*, 2019, pp. 5752–5758.
- [44] J. Geiping, Q. Garrido, P. Fernandez, A. Bar, H. Pirsiavash, Y. LeCun, and M. Goldblum, "A cookbook of self-supervised learning," *arXiv preprint arXiv:2304.12210*, 2023.
- [45] S. P. González-Sabbagh and A. Robles-Kelly, "A survey on underwater computer vision," *ACM Computing Surveys*, vol. 55, no. 13s, pp. 1–39, 2023.
- [46] K. Grorud-Colvert, J. Sullivan-Stack, C. Roberts, V. Constant, B. Horta e Costa, E. P. Pike, N. Kingston, D. Laffoley, E. Sala, J. Claudet *et al.*, "The mpa guide: A framework to achieve global goals for the ocean," *Science*, vol. 373, no. 6560, p. eabf0861, 2021.
- [47] G. Ha, "Coral species classification dataset," Available on Roboflow, 2022.
- [48] L. Haixin, Z. Ziqiang, M. Zeyu, and S.-K. Yeung, "Marinedet: Towards open-marine object detection," *arXiv preprint arXiv:2310.01931*, 2023.
- [49] B. S. Halpern, S. Walbridge, K. A. Selkoe, C. V. Kappel, F. Micheli, C. d'Agrosa, J. F. Bruno, K. S. Casey, C. Ebert, H. E. Fox *et al.*, "A global map of human impact on marine ecosystems," *science*, vol. 319, no. 5865, pp. 948–952, 2008.
- [50] Z. Hao, J. Qiu, H. Zhang, G. Ren, and C. Liu, "Umotma: Underwater multiple object tracking with memory aggregation," *Frontiers in Marine Science*, vol. 9, p. 1071618, 2022.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [52] J. Hong, M. Fulton, and J. Sattar, "Trashcan: A semantically-segmented dataset towards visual detection of marine debris," *CoRR*, vol. abs/2007.08097, 2020.
- [53] L. Hong, X. Wang, G. Zhang, and M. Zhao, "Usod10k: A new benchmark dataset for underwater salient object detection," *IEEE Transactions on Image Processing*, pp. 1–1, 2023.
- [54] Y. Hu, K. Wang, X. Zhao, H. Wang, and Y. Li, "Underwater image restoration based on convolutional neural network," in *Asian conference on machine learning*. PMLR, 2018, pp. 296–311.
- [55] S. Huang, K. Wang, H. Liu, J. Chen, and Y. Li, "Contrastive semi-supervised learning for underwater image restoration via reliable bank," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 18145–18155.
- [56] M. J. Islam, C. Edge, Y. Xiao, P. Luo, M. Mehtaz, C. Morse, S. S. Enan, and J. Sattar, "Semantic Segmentation of Underwater Imagery: Dataset and Benchmark," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE/RSJ, 2020.
- [57] M. J. Islam, Y. Xia, and J. Sattar, "Fast underwater image enhancement for improved visual perception," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3227–3234, 2020.
- [58] M. Jahidul Islam, R. Wang, and J. Sattar, "Svam: Saliency-guided visual attention modeling by autonomous underwater robots," *arXiv e-prints*, pp. arXiv–2011, 2020.
- [59] A. Jalal, A. Salman, A. Mian, M. Shortis, and F. Shafait, "Fish detection and species classification in underwater environments using deep learning with temporal information," *Ecological Informatics*, vol. 57, p. 101088, 2020.
- [60] S. Jennings and M. J. Kaiser, "The effects of fishing on marine ecosystems," in *Advances in marine biology*. Elsevier, 1998, vol. 34, pp. 201–352.
- [61] K. Katija, E. Orenstein, B. Schlining, L. Lundsten, K. Barnard, G. Sainz, O. Boulais, M. Cromwell, E. Butler, B. Woodward, and K. L. Bell, "FathomNet: A global image database for enabling artificial intelligence in the ocean," *Sci. Rep.*, vol. 12, no. 1, pp. 1–14, 2022.
- [62] J. Kay and M. Merrifield, "The fishnet open images database: A dataset for fish detection and fine-grained categorization in fisheries," *arXiv preprint arXiv:2106.09178*, 2021.
- [63] L. Ke, M. Danelljan, X. Li, Y.-W. Tai, C.-K. Tang, and F. Yu, "Mask transfiner for high-quality instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 4412–4421.
- [64] F. F. Khan, X. Li, A. J. Temple, and M. Elhoseiny, "Fishnet: A large-scale dataset and benchmark for fish recognition, detection, and functional trait prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 20496–20506.
- [65] M. U. Khattak, H. Rasheed, M. Maaz, S. Khan, and F. S. Khan, "Maple: Multi-modal prompt learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19113–19122.
- [66] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollar, and R. Girshick, "Segment anything," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 4015–4026.
- [67] A. Kirillov, Y. Wu, K. He, and R. Girshick, "Pointrend: Image segmentation as rendering," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [68] S.-H. Lee and M.-H. Oh, "Detection and tracking of underwater fish using the fair multi-object tracking model: A comparative analysis of yolov5s and dla-34 backbone models," *Applied Sciences*, vol. 14, no. 16, p. 6888, 2024.
- [69] J. Li, K. A. Skinner, R. M. Eustice, and M. Johnson-Roberson, "Watergan: Unsupervised generative network to enable real-time color correction of monocular underwater images," *IEEE Robotics and Automation letters*, vol. 3, no. 1, pp. 387–394, 2017.
- [70] J. Li, W. Xu, L. Deng, Y. Xiao, Z. Han, and H. Zheng, "Deep learning for visual recognition and detection of aquatic animals: A review," *Reviews in Aquaculture*, vol. 15, no. 2, pp. 409–433, 2023.
- [71] J. Li, D. Li, S. Savarese, and S. Hoi, "Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," in *International conference on machine learning*. PMLR, 2023, pp. 19730–19742.
- [72] X. Li, Y. Huang, Z. He, Y. Wang, H. Lu, and M.-H. Yang, "Citetracker: Correlating image and text for visual tracking," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 9974–9983.
- [73] Y. Li, F. Liang, L. Zhao, Y. Cui, W. Ouyang, J. Shao, F. Yu, and J. Yan, "Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm," *arXiv preprint arXiv:2110.05208*, 2021.
- [74] Y. Li, B. Wang, Y. Li, Z. Liu, W. Huo, Y. Li, and J. Cao, "Underwater object tracker: Uostrack for marine organism grasping of underwater vehicles," *Ocean Engineering*, vol. 285, p. 115449, 2023.
- [75] S. Lian, H. Li, R. Cong, S. Li, W. Zhang, and S. Kwong, "Watermask: Instance segmentation for underwater imagery," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 1305–1315.
- [76] D. Liang, W. Xu, and X. Bai, "An end-to-end transformer model for crowd localization," *European Conference on Computer Vision*, 2022.
- [77] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [78] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in

- context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [79] C. Liu, Z. Wang, S. Wang, T. Tang, Y. Tao, C. Yang, H. Li, X. Liu, and X. Fan, "A new dataset, poisson gan and aquanet for underwater object grabbing," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 5, pp. 2831–2844, 2022.
- [80] J. Liu, A. T. Becerra, J. F. Bienvenido-Barcena, X. Yang, Z. Zhao, and C. Zhou, "Cffii-vit: Enhanced vision transformer for the accurate classification of fish feeding intensity in aquaculture," *Journal of Marine Science and Engineering*, vol. 12, no. 7, p. 1132, 2024.
- [81] Y. Liu, D. An, Y. Ren, J. Zhao, C. Zhang, J. Cheng, J. Liu, and Y. Wei, "Dp-fishnet: Dual-path pyramid vision transformer-based underwater fish detection network," *Expert Systems with Applications*, vol. 238, p. 122018, 2024.
- [82] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 11976–11986.
- [83] I. Loshchilov, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017.
- [84] J. Lu, N. Li, S. Zhang, Z. Yu, H. Zheng, and B. Zheng, "Multi-scale adversarial network for underwater image restoration," *Optics & Laser Technology*, vol. 110, pp. 105–113, 2019.
- [85] M. Maniparambil, C. Vorster, D. Molloy, N. Murphy, K. McGuinness, and N. E. O'Connor, "Enhancing clip with gpt-4: Harnessing visual descriptions as prompts," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 262–271.
- [86] J. R. Merrick, "Australasian freshwater fish faunas: diversity, interrelationships, radiations and conservation," in *Evolution and biogeography of Australasian vertebrates*. Auscipub, 2006, pp. 195–224.
- [87] Z. Michel, "The rich biodiversity of ocean ecosystems," *National Geographic*, December 2022, accessed: November 9, 2024.
- [88] M. Modasshir, S. Rahman, O. Youngquist, and I. Rekleitis, "Coral identification and counting with an autonomous underwater vehicle," in *2018 IEEE international conference on robotics and biomimetics (ROBIO)*. IEEE, 2018, pp. 524–529.
- [89] L. Mohan, "Innovative approaches to coral conservation," *National Geographic*, June 2023, accessed: November 9, 2024.
- [90] Y. Mu, Q. Zhang, M. Hu, W. Wang, M. Ding, J. Jin, B. Wang, J. Dai, Y. Qiao, and P. Luo, "Embodiedgpt: Vision-language pre-training via embodied chain of thought," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [91] M. F. Naeem, M. G. Z. A. Khan, Y. Xian, M. Z. Afzal, D. Stricker, L. Van Gool, and F. Tombari, "I2mvformer: Large language model generated multi-view document supervision for zero-shot image classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 15169–15179.
- [92] M. Oquab, T. Dariset, T. Moutakanni, H. Vo, M. Szafrańiec, V. Khaldov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby *et al.*, "Dinov2: Learning robust visual features without supervision," 2023.
- [93] K. Panetta, L. Kezebou, V. Oludare, and S. Agaian, "Comprehensive underwater object tracking benchmark dataset and underwater image enhancement with gan," *IEEE Journal of Oceanic Engineering*, vol. 47, no. 1, pp. 59–75, 2022.
- [94] M. Pedersen, J. Bruslund Haurum, R. Gade, and T. B. Moeslund, "Detection of marine animals in a new underwater dataset with varying visibility," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- [95] L. Peng, C. Zhu, and L. Bian, "U-shape transformer for underwater image enhancement," *IEEE Transactions on Image Processing*, vol. 32, pp. 3066–3079, 2023.
- [96] C. Prathima, C. Silpa, A. Charitha, G. Harshitha, C. Sai Charan, and G. R. Sailendra, "Detecting and recognizing marine animals using advanced deep learning models," in *2024 International Conference on Expert Clouds and Applications (ICOECA)*, 2024, pp. 950–955.
- [97] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [98] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [99] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.
- [100] A. Saleh, I. H. Laradji, D. A. Konovalov, M. Bradley, D. Vazquez, and M. Sheaves, "A realistic fish-habitat dataset to evaluate algorithms for underwater visual analysis," *Scientific Reports*, vol. 10, no. 1, p. 14671, 2020.
- [101] A. Saleh, M. Sheaves, D. Jerry, and M. R. Azghadi, "Transformer-based self-supervised fish segmentation in underwater videos," *arXiv preprint arXiv:2206.05390*, 2022.
- [102] A. Saleh, M. Sheaves, and M. Rahimi Azghadi, "Computer vision and deep learning for fish classification in underwater habitats: A survey," *Fish and Fisheries*, vol. 23, no. 4, pp. 977–999, 2022.
- [103] F. Sammani, T. Mukherjee, and N. Deligiannis, "NLx-gpt: A model for natural language explanations in vision and vision-language tasks," in *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 8322–8332.
- [104] R. Schodde, *Zoological Catalogue of Australia: Aves (Columbidae to Coraciidae)*. CSIRO PUBLISHING, 1997, vol. 37.
- [105] X. Shao, H. Chen, K. Magson, J. Wang, J. Song, J. Chen, and J. Sasaki, "Deep learning for multi-label classification of coral conditions in the indo-pacific via underwater photogrammetry," *arXiv preprint arXiv:2403.05930*, 2024.
- [106] J. J. Shelley, A. Delaval, and M. C. Le Feuvre, "A revision of the grunter genus syncomistes (teleostei, terapontidae, syncomistes) with descriptions of seven new species from the kimberley region, northwestern australia," *Zootaxa*, vol. 4367, no. 1, pp. 1–103, 2017.
- [107] G. Si, Y. Xiao, B. Wei, L. B. Bullock, Y. Wang, and X. Wang, "Token-selective vision transformer for fine-grained image recognition of marine organisms," *Frontiers in Marine Science*, vol. 10, p. 1174347, 2023.
- [108] S. A. Siddiqui, A. Salman, M. I. Malik, F. Shafait, A. Mian, M. R. Shortis, and E. S. Harvey, "Automatic fish species classification in underwater videos: exploiting pre-trained deep neural network models to compensate for limited labelled data," *ICES Journal of Marine Science*, vol. 75, no. 1, pp. 374–389, 2018.
- [109] A. Singh, R. Hu, V. Goswami, G. Couairon, W. Galuba, M. Rohrbach, and D. Kiela, "Flava: A foundational language and vision alignment model," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15638–15650.
- [110] Q. Song, C. Wang, Z. Jiang, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, and Y. Wu, "Rethinking counting and localization in crowds: A purely point-based framework," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 3365–3374.
- [111] G. Sun, Z. An, Y. Liu, C. Liu, C. Sakaridis, D.-P. Fan, and L. Van Gool, "Indiscernible object counting in underwater scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13791–13801.
- [112] Z. Sun, Q. Ke, H. Rahmani, M. Bennamoun, G. Wang, and J. Liu, "Human action recognition from various data modalities: A review," *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 3, pp. 3200–3225, 2022.
- [113] S. Thornton, "Understanding coral bleaching and ocean health," *National Geographic*, July 2024, accessed: November 9, 2024.
- [114] S. Thornton and L. J. Richardson, "The vital role of coral reefs and the threats they face," *National Geographic*, July 2024, accessed: November 9, 2024.
- [115] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*. Springer, 2020, pp. 776–794.
- [116] O. Ulucan, D. Karakaya, and M. Turkan, "A large-scale dataset for fish segmentation and classification," in *2020 Innovations in Intelligent Systems and Applications Conference (ASYU)*. IEEE, 2020, pp. 1–5.
- [117] M. Ummar, F. A. Dharejo, B. Alawode, T. Mahbub, M. J. Piran, and S. Javed, "Window-based transformer generative adversarial network for autonomous underwater image enhancement," *Engineering Applications of Artificial Intelligence*, vol. 126, p. 107069, 2023.
- [118] R. Van Der Laan, W. N. Eschmeyer, and R. Fricke, "Family-group names of recent fishes," *Zootaxa*, vol. 3882, no. 1, pp. 1–230, 2014.
- [119] J. Veron, M. Stafford-Smith, E. Turak, and L. De Ventier, "Corals of the world. accessed 12 mar 2023, version 0.01," 2016.
- [120] J. Wan, Q. Wang, and A. B. Chan, "Kernel-based density map generation for dense object counting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 3, pp. 1357–1370, 2022.
- [121] N. Wang, T. Chen, S. Liu, R. Wang, H. R. Karimi, and Y. Lin, "Deep learning-based visual detection of marine organisms: A survey," *Neurocomputing*, vol. 532, pp. 1–32, 2023.
- [122] X. Wang, R. Zhang, T. Kong, L. Li, and C. Shen, "SOLOv2: Dynamic and fast instance segmentation," in *Adv. Neural Inf. Process. Syst.*, vol.

- 2020-December, 2020, Conference paper. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85107865426&partnerID=40&md5=e9e4f881a791ddbd81f564f74a345507>
- [123] Z. Wang, W. Liu, Y. Wang, and B. Liu, "Agcyclegan: Attention-guided cyclegan for single underwater image restoration," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 2779–2783.
 - [124] B. Xu, H. Liang, R. Liang, and P. Chen, "Locate globally, segment locally: A progressive architecture with knowledge review network for salient object detection," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 4, 2021, pp. 3004–3012.
 - [125] S. Xu, M. Zhang, W. Song, H. Mei, Q. He, and A. Liotta, "A systematic review and analysis of deep learning-based underwater object detection," *Neurocomputing*, vol. 527, pp. 204–232, 2023.
 - [126] W. Yan, Y. Zhang, P. Abbeel, and A. Srinivas, "Videogpt: Video generation using vq-vae and transformers," *arXiv preprint arXiv:2104.10157*, 2021.
 - [127] C.-Y. Yang, H.-W. Huang, Z. Jiang, H. Wang, F. Wallace, and J.-N. Hwang, "A density-guided temporal attention transformer for indiscernible object counting in underwater videos," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 5075–5079.
 - [128] L. Yang, Z. Xu, H. Zeng, N. Sun, B. Wu, W. Cheng, J. Bo, L. Li, Y. Dong, and S. He, "Fishdb: an integrated functional genomics database for fishes," *BMC Genomics*, vol. 21, 11 2020.
 - [129] Z. Yang, L. Li, K. Lin, J. Wang, C.-C. Lin, Z. Liu, and L. Wang, "The dawn of Imms: Preliminary explorations with gpt-4v (ision)," *arXiv preprint arXiv:2309.17421*, vol. 9, no. 1, p. 1, 2023.
 - [130] C.-H. Yeh, C.-H. Lin, L.-W. Kang, C.-H. Huang, M.-H. Lin, C.-Y. Chang, and C.-C. Wang, "Lightweight deep neural network for joint learning of underwater object detection and color conversion," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 11, pp. 6129–6143, 2021.
 - [131] M. Zand, H. Damirchi, A. Farley, M. Molahasani, M. Greenspan, and A. Etemad, "Multiscale crowd counting and localization by multitask point supervision," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 1820–1824.
 - [132] C. Zhang, R. Cong, Q. Lin, L. Ma, F. Li, Y. Zhao, and S. Kwong, "Cross-modality discrepant interaction network for rgb-d salient object detection," in *Proceedings of the 29th ACM international conference on multimedia*, 2021, pp. 2094–2102.
 - [133] C. Zhang, L. Liu, G. Huang, H. Wen, X. Zhou, and Y. Wang, "Webuot-1m: Advancing deep underwater object tracking with a million-scale benchmark," *arXiv preprint arXiv:2405.19818*, 2024.
 - [134] J. Zhang, J. Huang, S. Jin, and S. Lu, "Vision-language models for vision tasks: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
 - [135] P. Zhang, T. Yan, Y. Liu, and H. Lu, "Fantastic animals and where to find them: Segment any marine animal with dual sam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 2578–2587.
 - [136] W. Zhang, G. Chen, P. Zhuang, W. Zhao, and L. Zhou, "Catnet: Cascaded attention transformer network for marine species image classification," *Expert Systems with Applications*, p. 124932, 2024.
 - [137] Z. Zhao, C. Xia, C. Xie, and J. Li, "Complementary trilateral decoder for fast and accurate salient object detection," in *Proceedings of the 29th acm international conference on multimedia*, 2021, pp. 4967–4975.
 - [138] Z. Zheng, Y. Chen, J. Zhang, T.-A. Vu, H. Zeng, Y. H. W. Tim, and S.-K. Yeung, "Exploring boundary of gpt-4v on marine analysis: A preliminary case study," *arXiv preprint arXiv:2401.02147*, 2024.
 - [139] Z. Zheng, H. Liang, B.-S. Hua, Y. H. Wong, P. Ang, A. P. Y. Chui, and S.-K. Yeung, "Coralscop: Segment any coral image on this planet," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 28170–28180.
 - [140] Z. Zheng, J. Zhang, T.-A. Vu, S. Diao, Y. H. W. Tim, and S.-K. Yeung, "Marinegpt: Unlocking secrets of ocean to the public," *arXiv preprint arXiv:2310.13596*, 2023.
 - [141] J. Zhong, M. Li, H. Zhang, and J. Qin, "Combining photogrammetric computer vision and semantic segmentation for fine-grained understanding of coral reef growth under climate change," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 186–195.
 - [142] Y. Zhong, J. Yang, P. Zhang, C. Li, N. Codella, L. H. Li, L. Zhou, X. Dai, L. Yuan, Y. Li *et al.*, "Regionclip: Region-based language-image pretraining," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16793–16803.
 - [143] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.
 - [144] Q. Zhou, S. Wang, J. Liu, X. Hu, Y. Liu, Y. He, X. He, and X. Wu, "Geological evolution of offshore pollution and its long-term potential impacts on marine ecosystems," *Geoscience Frontiers*, vol. 13, no. 5, p. 101427, 2022.
 - [145] C. Zhu, R. H. Byrd, P. Lu, and J. Nocedal, "Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale bound-constrained optimization," *ACM Transactions on mathematical software (TOMS)*, vol. 23, no. 4, pp. 550–560, 1997.
 - [146] P. Zhu, H. Wang, and V. Saligrama, "Don't even look once: Synthesizing features for zero-shot detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11693–11702.
 - [147] P. Zhuang, Y. Wang, and Y. Qiao, "Wildfish++: A comprehensive fish benchmark for multimedia research," *IEEE Transactions on Multimedia*, vol. 23, pp. 3603–3617, 2021.
 - [148] Z. Ziqiang, X. Yaofeng, L. Haixin, Y. Zhibin, and S.-K. Yeung, "Coralvos: Dataset and benchmark for coral video segmentation," *arXiv preprint arXiv:2310.01946*, 2023.
 - [149] Z. Ziqiang, C. Yiwe, Z. Huimin, V. Tuan-Anh, H. Bin-Son, and Y. Sai-Kit, "Marineinst: A foundation model for marine image analysis with instance visual description," *European Conference on Computer Vision (ECCV)*, 2024.



Basit Alawode received his BSc degree in electrical engineering from the Obafemi Awolowo University, Nigeria, in 2013. He further obtained his MSc degree in electrical engineering at the King Fahd University of Petroleum and Minerals, Saudi Arabia in 2020. He is currently pursuing his Ph.D. degree in computer science and engineering degree at Khalifa University of Science and Technology, UAE. His research interests include visual object tracking and computer vision.



Iyyakutti Iyappan Ganapathi is currently a postdoctoral fellow at Department of Electrical Engineering and Computer Science, Khalifa University, Abu Dhabi, UAE. He previously worked as an Assistant Professor at Woosong University in South Korea. He earned his PhD degree from the Indian Institute of Technology Indore, India. 3D image processing, biometrics, computer vision, and machine learning are among his research interests.



Sajid Javed is a faculty member at Khalifa University (KU), UAE. Prior to that, he was a research fellow at KU from 2019 to 2021 and at the University of Warwick, UK, from 2017–2018. He received his B.Sc. degree in computer science from the University of Hertfordshire, UK, in 2010. He completed his combined Master's and Ph.D. degrees in computer science from Kyungpook National University, Republic of Korea, in 2017.



Arif Mahmood is currently a professor at the Department of Computer Science at ITU, Pakistan, and the director of the computer vision lab. He received his M.Sc. and Ph.D. degree in computer science from LUMS, Pakistan, in 2003 and 2011. He has also worked as a Research Assistant Professor with the School of Mathematics and Statistics (SMS), University of Western Australia (UWA). His research interests are face recognition, object classification, human-object interaction detection, and abnormal event detection.



Naoufel Werghi is a Professor at the department of computer science in Khalifa University for Science and Technology, UAE. He received his Habilitation and PhD in Computer Vision from the University of Strasbourg. His main research area is 2D/3D image analysis and interpretation, where he has been leading several funded projects related to biometrics, medical imaging, remote sensing, and intelligent systems.

Supplementary Material AquaticCLIP: A Vision-Language Foundation Model for Underwater Scene Analysis

TABLE OF CONTENTS

- 1) Construction, Curation, and Cleaning of a 2M Aquatic Image-Text Paired Dataset (Sec. VI).
- 2) Additional Training and Implementation Details (Sec. VII).
- 3) More Ablation Studies (Sec. VIII).
- 4) MRegionCLIP Pre-training for Instance Detection (Sec. IX).
- 5) Unsupervised Generation of Image and Instance-Level Descriptions (Sec. X).
- 6) Underwater Datasets and Tasks (Sec. XI).
- 7) SOTA Methods for Comparison (Sec. XII).
- 8) SOTA Methods Training Details (Sec. XIII).
- 9) Evaluation Metrics (Sec. XIV).
- 10) Zero-shot Inference (Sec. XV).
- 11) Zero-shot Cross-Modal Retrieval Results (Sec. XVI).
- 12) Underwater Object Detection and Classification Results (Sec. XVII).
- 13) Underwater Scene Segmentation Results (Sec. XVIII).
- 14) Results of Object Counting in Underwater Scenes (Sec. XIX).
- 15) Supervised AquaticCLIP Model: Linear Probe Evaluations (Sec. XX).
- 16) AquaticVision: Vision-Only Model Pre-training Details (Sec. XXI).
- 17) Visual Results (Sec. XXII).
- 18) Why AquaticCLIP Performance is Better? (Sec. XXIII).

VI. CONSTRUCTION, CURATION, AND CLEANING OF A 2M AQUATIC IMAGE-TEXT PAIRED DATASET

Our aquatic dataset consists of a diverse collection of 2 million image-text pairs collected from various online and heterogeneous sources. These include YouTube and Netflix videos, a variety of textbooks on marine biology and oceanography, content on underwater species like fish and sharks, marine articles, as well as repositories like Coral of the World [119], Fishes of Australia [86], [104], [106], [118], Marine Twitter, and National Geographic (NatGeo) [32], [36], [40], [87], [89], [113], [114]. In total, the dataset contains 20.3 million instances, with at least one instance per image and an average of 10.3 instances per image. To supplement freely available resources, we allocated \$6550 USD toward subscriptions, PDFs, CDs of marine biology books, and NatGeo magazines. Figs. 4-5 show some of the exemplar image-text pairs from our 2M dataset. Fig. 6 illustrates the dataset construction pipeline, which includes two key steps: the collection of image-text pairs from multiple sources, followed by meticulous cleaning and filtering to ensure dataset quality.

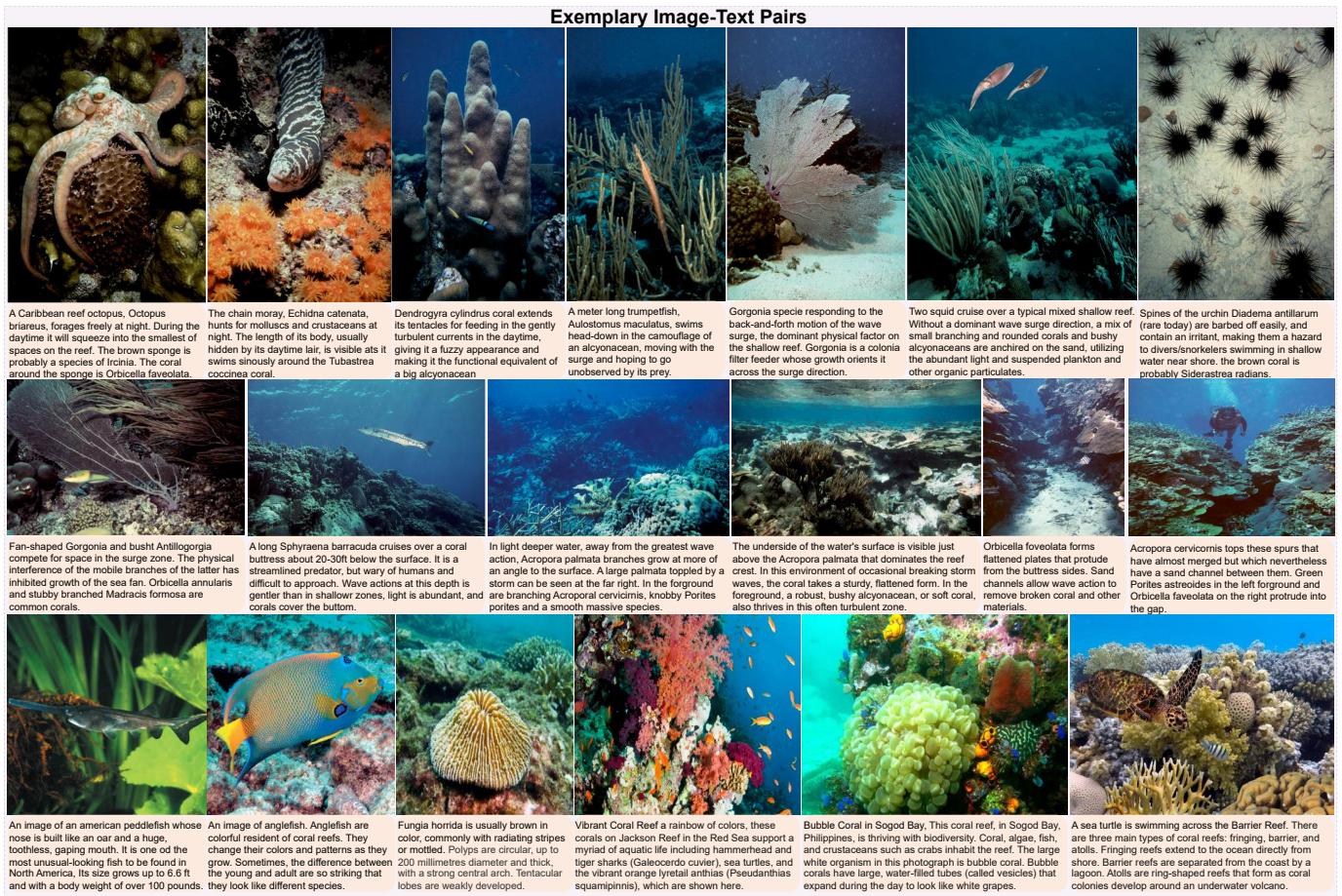


Fig. 4: Exemplary image-text pairs from our 2 million aquatic image-text paired dataset.



Fig. 5: Exemplary image-text pairs from our 2 Million aquatic image-text paired dataset.

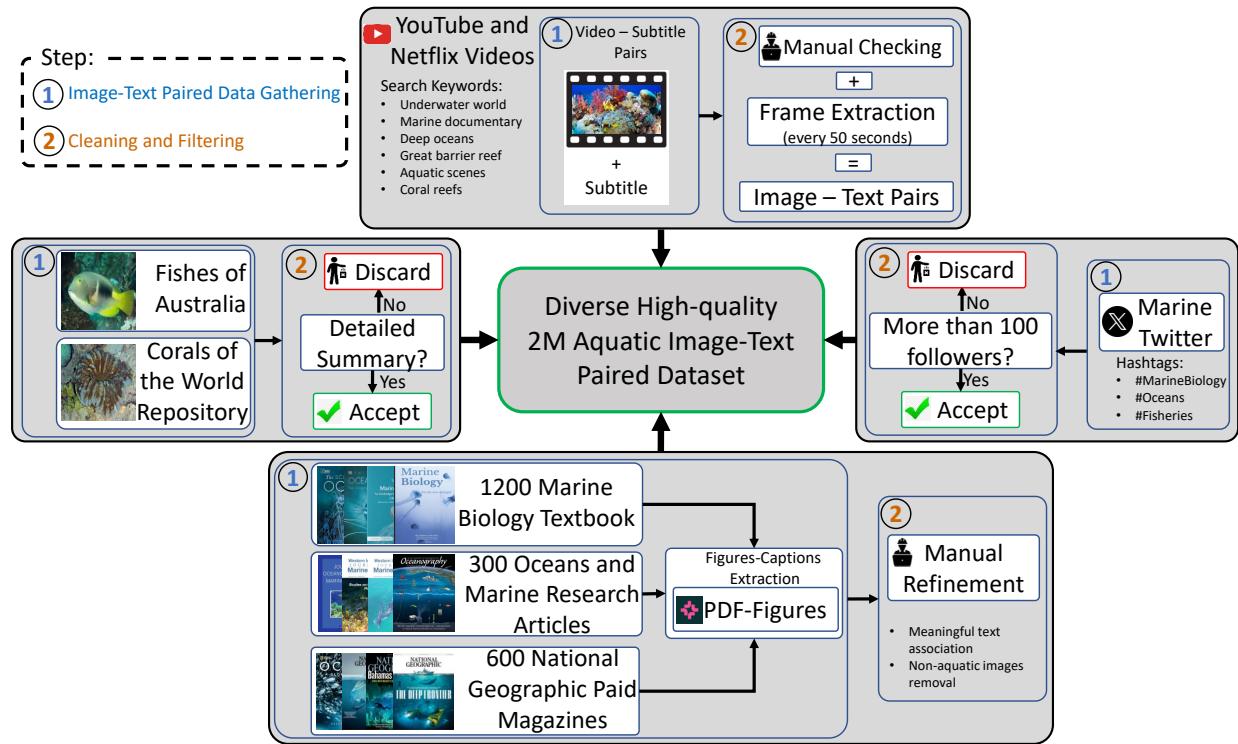


Fig. 6: Workflow for constructing and cleaning a diverse, high-quality 2 million aquatic image-text paired dataset for Aquatic-CLIP model pre-training. This process involves two main steps: (1) Image-Text Paired Data Gathering from sources such as YouTube and Netflix videos, marine biology textbooks, online repositories, and social media, and (2) Cleaning and Filtering through manual checks, frame extraction, caption refinement, and content validation to ensure dataset accuracy and relevance..

A. YouTube Videos

We reviewed several hundred documentary videos on underwater, marine, and ocean topics, using keywords such as “underwater world”, “underwater paradise”, “Marine Life and Ocean Animals”, “Our Planet”, “Wonders of Oceans”, “underwater killers”, “sea creatures”, “ocean animals for kids”, “kingdoms of marine life”, “Reef Life of the Andaman”, “The Big Blue Ocean stories”, “marine documentary”, “deep oceans”, “great barrier reef”, “aquatic scenes”, “coral reefs”, “Gorgeous Underwater World of the Philippines”, “Mysterious hunters of the deep sea”, “Dark side of the ocean”, “Beautiful Coral Reefs”, “Mysterious of the Twilight Zone”.

B. Netflix Documentaries

Using a paid subscription, we explored hundreds of hours of documentary videos, using the following titles: “My Octopus Teacher”, “Last Breath”, “Wonders of the Reef”, “Seaspiracy”, “Mission blue”, “The deepest breath”, “Our planet”, “chasing coral”, “Microworlds reef”, “Secret world of sound”, “Night of earth”, “Disney Nature Oceans”, “A Plastic Ocean”, “Our Oceans”.

Subtitles from these videos were used to create aligned image-text pairs, which were manually reviewed and refined. We extracted unique frames every 50 seconds, which presented challenges like low visibility, motion blur, background clutter,

and color distortions. In total, we captured 0.8 million image-text pairs from YouTube videos and 0.5 million pairs from Netflix documentaries.

C. Textbooks on Marine Biology, Oceanography, Underwater Creatures and Species, Fish, and Sharks

We used a total of 1,200 different textbooks on marine biology, ocean science, underwater creatures and species, fish, sea life, and coral reefs, yielding 0.4 million image-text pairs. Below are 50 example books:

1. Exploring the World of Aquatic Life by John Dawes.
2. Marine Biology Basics.
3. New Species Described in Corals of the World by Veron.
4. Introduction to Marine Biology, 3rd Edition by George.
5. Marine Biology: A Very Short Introduction by Philip.
6. The soul of an octopus: A surprising exploration into the wonder of consciousness by Montgomery.
7. The World Beneath: The Life and Times of Unknown Sea Creatures and Coral Reefs by Richard.
8. Becoming a Marine Biologist by Morell.
9. Handbook of Whales, Dolphins, and Porpoises of the World by Mark.
10. Marine Biology: An Ecological Approach by James.
11. Citizens of the Sea: Wondrous Creatures from the Census of Marine Life by Nancy.
12. Deep: Freediving, Renegade science, and what the ocean tells us about ourselves by Nestor.
13. Marine Biology For The Non-Biologist by Andrew.
14. Marine Biology: Comparative Ecology of Planet Ocean by Roberto.
15. Marine Ecology by Michel.
16. Oceanography and marine biology: an introduction to marine science by David.
17. Watching giants: the

secret lives of whales by Kelsey. 18. *The ocean of life: The fate of man and the sea by Callum.* 19. *Sink Like Fish: How Sound Rules Life Underwater by Amorina.* 20. *The Marine World by Frances.* 21. *Marine Biology: Function, Biodiversity, Ecology by Jeffrey S. Levinton.* 22. *Introduction to the Biology of Marine Life by John Morrissey and James L. Sumich.* 23. *Marine Biology by Peter H. Raven and George B. Johnson.* 24. *Biology of the Marine Environment by John D. H. Connell.* 25. *Marine Biology: A Very Short Introduction by Philip V. Wells.* 26. *Marine Biology for the Non-Biologist by Anne E. McCarthy.* 27. *Fundamentals of Marine Biology by Paul F. McCarthy.* 28. *Marine Conservation Biology by H. A. J. W. Davis.* 29. *Marine Biotechnology by David L. K. A. V. Hunter.* 30. *Tropical Marine Ecology by R. G. K. Smiley.* 31. *Marine Biogeochemistry: A Multidisciplinary Approach by K. J. C. Browne.* 32. *Fisheries Biology, Assessment and Management by A. A. G. J. K. Smith.* 33. *Marine Biotechnology: Methods and Protocols by C. B. H. R. S. Farago.* 34. *Aquatic Ecosystems: Trends and Global Perspectives by A. C. W. F. R. B. Stokes.* 35. *Marine Ecotoxicology by T. B. R. R. O. G. C. Chikaram.* 36. *Marine Biogeography: An Overview by M. C. E. D. E. K. P. Chen.* 37. *Ocean Biogeochemistry by M. R. M. B. D. B. F. N. Andrews.* 38. *Marine Ecology: A Comprehensive Study of the Ecology of Marine Organisms by M. J. E. S. McMahon.* 39. *Aquaculture: Principles and Practices by D. S. S. H. L. P. T. J. S. M. G. Thomas.* 40. *Marine Resource Management: A Global Perspective by L. J. T. S. B. W. Smith.* 41. *Marine Policy: A Comprehensive Overview by R. B. M. W. G. H. J. K. M. W. D. J. H. Mark.* 42. *Marine Ecosystem Management by D. J. M. J. K. R. K. M. S. T. J. Stokes.* 43. *The Diversity of Fishes: Biology, Evolution, and Ecology.* 44. *Marine Invertebrates of the Pacific Northwest.* 45. *Coral Reef Fishes: Dynamics and Diversity in a Complex Ecosystem.* 46. *Introduction to Marine Biology.* 47. *Fishes of the Gulf of Maine.* 48. *Marine Mammals: Evolutionary Biology.* 49. *Marine Algae of the Pacific Northwest.* 50. *The Ecology of Fishes on Coral Reefs.*

D. Oceans and Marine Research Articles

We gathered 178K image-text pairs from 300 marine science articles obtained through Google.

E. National Geographic (NatGeo) Subscribed Magazines

We collected 100K image-text pairs from 600 subscribed Nat-Geo magazines. Some sample magazine titles include “*Corals Reef*”, “*underwater exploration*”, “*national geographic explorer*”, “*national geographic sea change*”, “*Oceanography*”, “*Adventure*”, “*Science of the Sea*”, “*Magazine: The Ocean*”, “*The Wonders of the Sea*”, “*Oceans and the Environment*”, “*Underwater World*”, “*Marine Conservation*”, “*The Blue Planet*”, “*Exploring the Deep Sea*”, “*Life Underwater*”.

Figures and captions were extracted using PDF-Figures 2.0 tool [30], and we manually refined the data to ensure that the selected images had meaningful text and captions. Non-aquatic images were discarded.

F. Corals of the World and Fishes of Australia

These repositories are publicly available. We extracted 15K image-text pairs from these resources and manually verified that the images depicted aquatic scenes and had detailed, meaningful textual descriptions.

G. Marine Twitter

We used Twitter to search for relevant content using hashtags such as #MarineBiology, #Oceans, and #Fisheries, considering only channels with over 100 followers. We collected 7K image-text pairs from this source after a thorough cleaning and filtering process.

VII. ADDITIONAL TRAINING AND IMPLEMENTATION DETAILS

The AquaticCLIP model architecture includes a frozen domain-specific caption generator (MarineGPT [140]) to produce unsupervised additional textual descriptions, the CLIP [97] image encoder with ViT-B/16-224 [38] for extracting patch-level embeddings, and a transformer-based text encoder with a 76-token maximum sequence length [98] for textual embeddings. We fine-tuned four components: the image encoder, text encoder, prompt-guided vision encoder, and vision-guided text encoder, utilizing cross-modal contrastive loss as described in Section 3.6. The Adam optimizer [83] was used with an initial learning rate of 1×10^{-4} and weight decay of 1×10^{-5} . The model was trained over 80 epochs on four A100 GPUs, with a batch size of 512. We set 20 prompts for the model. During pre-training, all images were resized to 512×512 , with larger images resized on the short side to 512 and center-cropped. Data augmentation, including horizontal and vertical flips, was applied to both images and captions. A linear projection head mapped the text and image embeddings into a 512-dimensional latent space for alignment. Image and text representations were aligned using the cross-modal contrastive loss (Section 3.6), and the model was implemented with PyTorch.

In a second experiment, we fine-tuned only the image and text encoders, Φ_v and Φ_t , using a setup similar to CLIP [97] with our 2M image-text paired dataset. This configuration included a batch size, a weight decay of 0.2, a temperature setting of 0.07, a peak learning rate of 1×10^{-4} , the AdamW optimizer with an initial learning rate of 5×10^{-6} , and a cosine decay scheduler.

VIII. MORE ABLATION STUDIES

Figs. 13-19 show the architectural ablation figures highlighting the importance of the proposed Prompt-guided Vision Encoder (PGVE) and Vision-guided Text Encoder (VGTE). Please refer to Table 1 of the main manuscript.

- 1) **Comparison of Heterogeneous Textual Sources (Table V):** We evaluated AquaticCLIP with text descriptions generated by MarineGPT (MGPT) [140], GPT4V [129], and BLIP2 [71]. MGPT achieved the best results due to its specific pre-training in the aquatic domain, while GPT4V and BLIP2 were trained on general-domain data.

TABLE V: **Ablation study** comparing ground truth (GT) text with generated text from MarineGPT (MGPT), GPT4V, and BLIP2 for zero-shot classification (F_1 scores). The combination of GT and MGPT yields the best performance across all datasets, with top scores on MAI (0.871), SAI (0.923), and CC (0.953). Generated text from GPT4V and BLIP2 shows lower performance, highlighting the advantage of using domain-specific MGPT with GT..

Variants	GT	GT+MGPT	MGPT	GPT4V	BLIP2
MAI	0.861	0.871	0.844	0.811	0.801
SAI	0.902	0.923	0.883	0.834	0.818
FishNet	0.821	0.842	0.815	0.791	0.788
F NOI	0.785	0.801	0.769	0.743	0.727
LSF	0.917	0.934	0.905	0.881	0.880
CC	0.936	0.953	0.917	0.876	0.853

TABLE VI: **Ablation study** on AquaticCLIP, detailing the impact of varying the number of learnable prompts n_r in the Prompt-Guided Vision Encoder (PGVE) and the percentage of top-p% keywords. F_1 scores for zero-shot classification on the FishNet dataset indicate optimal performance at 20 prompts and 30% keyword retention, both achieving an F_1 score of 0.842. Changes in these parameters lead to minor performance variations.

Prompts (n_r)	5	10	15	20	25	30
AquaticCLIP	0.822	0.835	0.839	0.842	0.840	0.841
Top-p% keywords	90	70	50	30	20	10
AquaticCLIP	0.791	0.811	0.830	0.837	0.842	0.826

Combining ground truth with MGPT-generated texts further improved performance, as the two sources complemented each other in providing more accurate textual descriptions.

- 2) **Effect of the Number of Learnable Prompts on Performance (Table VI):** We varied the number of learnable prompts n_r in the Prompt-Guided Vision Encoder (PGVE) from 5 to 30, as shown in VI. Performance increased with more prompts, peaking at 0.842 for $n_r = 20$ on the FishNet dataset. Fewer prompts led to incomplete visual feature representation, while more than 20 prompts added redundancy without further performance gains.
- 3) **Performance Variation with Different Top-p% Keywords (Table VI):** In the image-text caption cleaning module, we experimented with retaining different percentages of top-matching keywords, ranging from 90% to 10%, as shown in Table VI. Optimal performance occurred when 20% of the top keywords were retained. Higher percentages likely included noisy keywords, reducing performance, while lower percentages discarded too many vital keywords, leading to a performance decrease.

IX. MREGIONCLIP PRE-TRAINING FOR INSTANCE DETECTION

The instance detection method, MRegionCLIP, is pre-trained on publicly available datasets (see Table VII) and then applied as a frozen instance detector to our 2M aquatic dataset.

We fine-tuned an existing object detector, RegionCLIP [142], using the frozen CLIP text encoder and ResNet50 as the vision backbone. Similar to methods like RegionCLIP [142], MarineDET [48], and MarineInst20M [149], the COCO Cap-

tion dataset was used for contrastive learning. For marine-specific pre-training, we fine-tuned on publicly available marine datasets, which contain 228,363 images and 664,411 instances. A powerful region proposal network (RPN) and a frozen image encoder from RegionCLIP were used. This fine-tuned model, named Marine RegionCLIP (MRegionCLIP), followed the same experimental protocols as MarineDET [48] and MarineInst20M [149]. The initial learning rate was set to 5×10^{-2} , then gradually reduced to 5×10^{-4} and 5×10^{-6} as needed. Fine-tuning was performed over 9,000 iterations with a batch size of 256 on four A100 GPUs. Quantitative object detection results for MRegionCLIP are provided in Table 10 of the main paper. MRegionCLIP was then utilized in zero-shot settings to detect instances in our 2M image-text paired dataset, and these detected instances were used for instance-level caption generation by MarineGPT [140].

X. UNSUPERVISED GENERATION OF IMAGE AND INSTANCE-LEVEL DESCRIPTIONS

Using our pre-trained object detector, MRegionCLIP, we detect objects as instances in each image within our 2M dataset. Each detected object instance is input into the MarineGPT model for caption generation. Specifically, we used the following prompt template: “The image is < image >. Describe the object in this image:”, where < image > is the image token. At the image level, each complete image I_i is fed into MarineGPT to generate corresponding textual descriptions. Fig. 7 illustrate the processes of image-level and instance-level caption generation, utilizing MRegionCLIP for object instance detection and MarineGPT for caption generation [140].

XI. UNDERWATER DATASETS AND TASKS

A. Zero-shot Marine Species Classification Datasets

For zero-shot classification tasks, we used the following datasets:

- 1) **Marine Animal Images Dataset [7]:** This dataset contains 806 images across nine different categories of sea animals, including Fish, Goldfish, Harbor Seal, Jellyfish, Lobster, Oyster, Sea Turtle, Squid, and Starfish. The images vary in size, lighting conditions, backgrounds, and camera angles, with an average resolution of 1024×768 . The dataset is split into 621 training images and 185 test images.
- 2) **Sea Animals Images Dataset [8]:** This dataset includes 13,711 images covering 23 different sea animals, such as Clams, Corals, Crabs, Dolphins, Eels, Fish, Jellyfish, Lobsters, Nudibranchs, Octopuses, Otters, Penguins, Puffers, Sea Rays, Sea Urchins, Seahorses, Seals, Sharks, Shrimps, Squids, Starfish, Turtles, and Whales, with an average resolution of 300×225 .

B. Zero-shot Fine-Grained Classification Datasets

For these experiments, we used the following datasets:

- 1) **FishNet [64]:** This large-scale, diverse dataset contains 94,532 images of fishes from 17,357 aquatic species, organized according to biological taxonomy (order, family,

TABLE VII: Publicly available underwater datasets used for pre-training MRegionCLIP in instance detection. Each dataset includes details on diversity, original task, number of images, instances, and average instances per image. The combined dataset totals 228,363 images and 664,678 instances, averaging 2.91 instances per image.

Dataset	Diversity	Original task and Motivation	Images	Instances	Average
Online images [4], [5], [9]	High	Image collection (human labeled)	35,175	194,010	5.52
Aquarium [3]	Medium	Underwater object detection	632	4,182	6.62
FathomNet [61]	High	Underwater and deep-sea object detection	69,909	121,329	1.74
FLOW [27]	Medium	Litter detection	1,825	3,850	2.39
MarineDET [48]	High	Open-marine object detection	22,679	39,243	1.73
MarineFouling [28]	Low	Biological fouling detection	221	508	2.30
OZFish [12]	Medium	Underwater fish detection	6,235	38,875	6.23
TACO [43]	Medium	Litter detection	1,109	2,656	2.39
TrashCAN [52]	Medium	Underwater trash detection	6,465	9,855	1.52
UDD/DUO [79]	Medium	Underwater object detection	2,170	13,090	6.03
Underwater_Garbage [10]	Medium	Underwater garbage detection	4,542	9,386	2.07
EOL [11]	High	Species identification	23,141	80,128	3.46
FishDB [128]	Medium	Fish species identification	9,905	18,914	1.91
HK-Reef-Fish [6]	Low	Fish identification	729	1,985	2.72
ImageNet_Sub [33]	Low	Scene classification	3,987	7,175	1.78
Oceanic_Life [14]	High	Collection of marine life imagery	5,029	20,811	4.14
Reef-Life-Survey [1]	High	Marine creature identification	7,075	12,502	1.77
Reeflex [2]	High	Marine creature identification	15,088	61,656	4.09
Sea Animals [8]	Medium	Sea animal classification	3,080	7,448	2.42
WildFish++ [147]	High	Fine-grained fish classification	9,367	17,075	1.82
Total	High	Instance detecting pre-training dataset	228,363	664,678	2.91

genus, and species). The dataset covers 83 orders and 463 families, with an average resolution of 542×372 . We report classification performance at both the family and order levels.

- 2) **FishNet Open Images dataset [62]:** This dataset consists of 86,029 images across 34 object classes, representing a large and diverse public dataset of fisheries with an average resolution of 1261×722 . It poses several challenges, including visual similarity between species, imbalanced class distributions, harsh weather conditions, and chaotic crew activities.
- 3) **Large Scale Fish Dataset [116]:** This dataset includes images of 9 different types of seafood collected from a supermarket: gilt-head bream, red sea bream, sea bass, red mullet, horse mackerel, black sea sprat, striped red mullet, trout, and shrimp. Images have an average resolution of 1688×1267 . The dataset contains a total of 18,000 images, with 2,000 samples per class after data augmentation

C. Zero-shot Coral Species Classification Datasets

This experiment was conducted using the following two datasets:

- 1) **Coral Species Classification (CSC) Dataset [47]:** This dataset contains 16 coral species classes with a total of 202 images, averaging a resolution of 257×192 . The classes are Acanthastrea-echinata, Acropora-millepora, Astreopora-myriophthalma, Coscinaraea-columna, Cyphastrea-microphthalma, Diploastrea-heliopora, Favites-pentagona, Goniopora-lobata, Montipora-stellata, Pavona-cactus, Plesiastrea-versipora, Pocillopora-

damicornis, Porites-lobata, Psammocora-contigua, Stylophora-pistillata, and Turbinaria-peltata.

- 2) **Corals Classification (CC) Dataset [35]:** This dataset consists of 9,292 images categorized into two classes: Healthy and Bleached Corals. The images have an average resolution of 294×231 and are divided into training, testing, and validation splits in an 80-10-10 ratio, resulting in 7,384 training images, 923 testing images, and 985 validation images.

D. Downstream Tasks and Datasets

We evaluated our AquaticCLIP model on several downstream analysis tasks, including salient object segmentation, instance segmentation, semantic segmentation, object detection and classification, and object counting.

For supervised salient object segmentation on underwater images, we used the USOD10K dataset [53]. For instance segmentation of marine images, we employed the Underwater Image Instance Segmentation (UIIS) dataset [75], and for semantic segmentation, we used the SUIM dataset [56]. Fine-tuned underwater object detection and classification tasks were conducted with four datasets: FishNet [64], DeepFish [100], URPC [13], and Brackish [94]. For marine object counting in biodiversity studies, we used the IOCFish5K dataset [111].

1) Instance, Salient Object, and Semantic Segmentation

Datasets: The fine-tuned segmentation experiments were performed on the following datasets:

- a) **Underwater Image Instance Segmentation (UIIS) [75]:** This dataset contains 4,628 images across 7 categories with pixel-level annotations for underwater instance segmentation. Categories include fish, coral

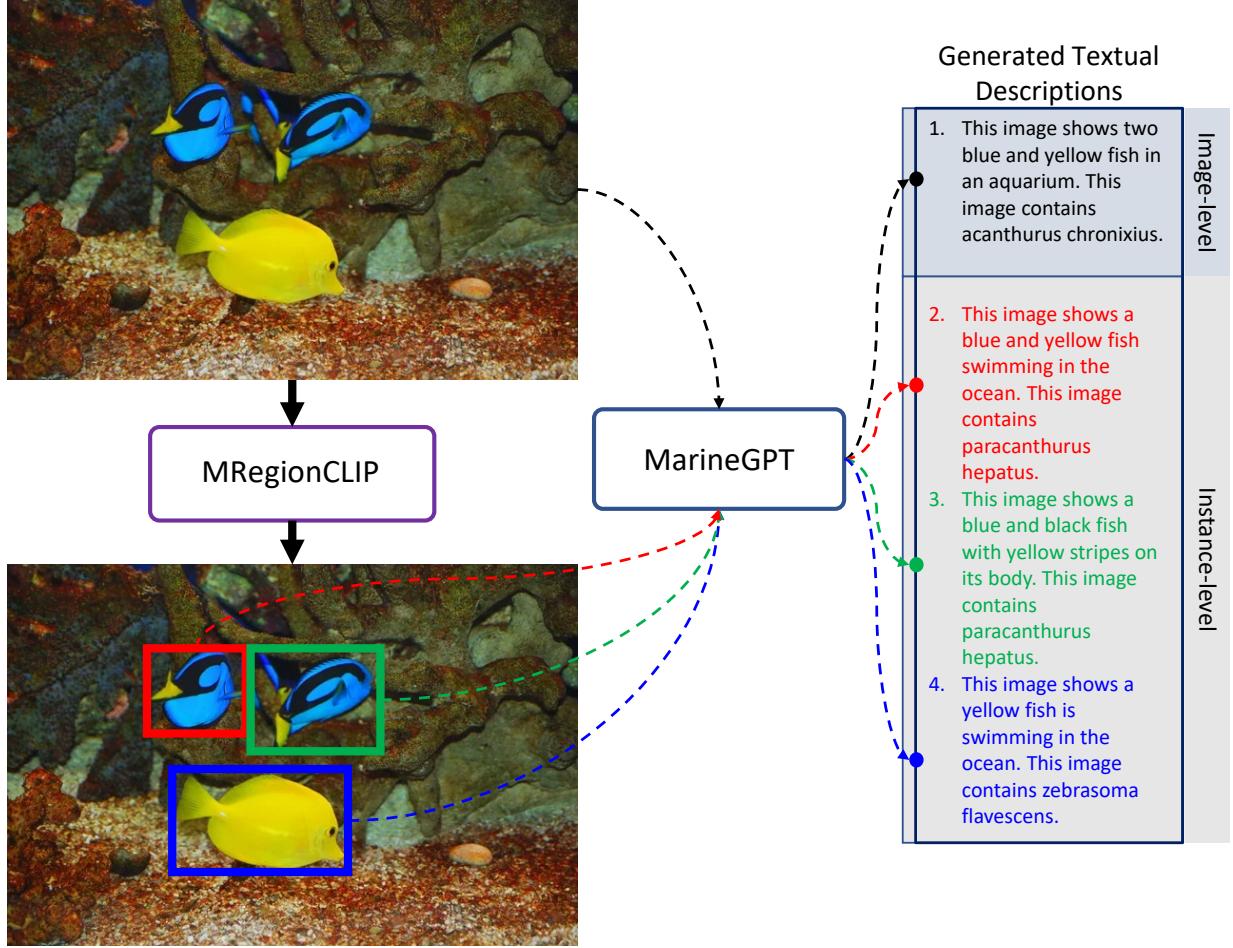


Fig. 7: Process of unsupervised image-and instance-level caption generation on a sample image using MRegionCLIP for object detection and MarineGPT [140] for generating descriptions. The detected instances (highlighted in colored boxes) are individually processed to produce detailed textual descriptions, with separate captions generated for the entire image (image-level) and each detected object (instance-level).

- reefs, aquatic plants, and wrecks/ruins, which are crucial for marine exploration and ecological studies. The images vary in resolution from 240×320 to 720×1280 .
- b) **USOD10K** [53]: This dataset is used for underwater salient object segmentation and contains 10,255 underwater images with a resolution of 640×480 pixels, covering 70 categories of salient objects across 12 underwater scenes. The dataset is divided into 7,178 training images, 2,051 validation images, and 1,026 testing images.
 - c) **Semantic Segmentation of Underwater IMagery (SUIM)** [56]: This dataset includes 1,635 image-mask pairs across 8 object categories for segmentation, including human divers, aquatic plants and seagrass, wrecks and ruins, robots, reefs and invertebrates, fish and vertebrates, sea-floor and rocks, and background. The masks are available in binary and combined RGB formats. The dataset is split into 1,525 training and 110 validation samples, with an average image resolution of 294×231 pixels.
 - 2) **Object Detection and Classification Datasets:** The object detection and classification tasks were performed on the following datasets:

- a) **FishNet** [64]: Previously discussed in the zero-shot classification task section.
- b) **DeepFish Dataset** [100]: Contains 39,766 images for tasks such as localization, counting, segmentation, and classification, collected from 20 habitats in tropical Australian marine environments. The images are captured in full HD resolution (1920×1080) from digital cameras. The dataset is split as follows: 19,883 images for training, 7,953 for validation, and 11,930 for testing in fish classification; 1,600, 640, and 960 for fish detection; and 310, 124, and 186 for fish segmentation.
- c) **Brackish Dataset** [94]: Contains high-resolution (1920×1080) images with bounding box annotations for 6 categories (big fish, small fish, starfish, shrimps, jellyfish, and crabs) captured in a brackish strait with variable visibility. It consists of 14,518 images, divided into training, validation, and test sets with an 80%:10%:10% split.
- d) **URPC Dataset** [13]: This dataset contains 6,626 high-resolution real underwater images with an average resolution of $2,388 \times 1,345$ pixels. The images feature four types of underwater objects: holothurian, echinus, scallop, and starfish. Identifying and distinguishing

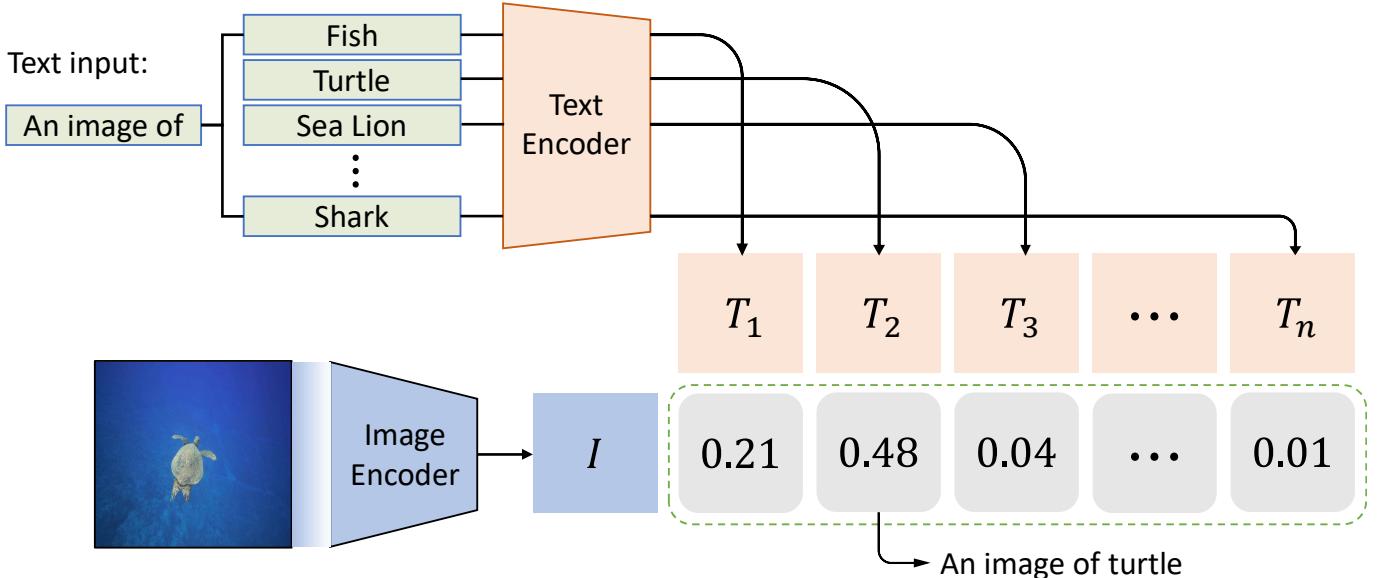


Fig. 8: Zero-shot inference using the proposed AquaticCLIP model. The process involves feeding a text input (e.g., “An image of”) with various marine categories (e.g., Fish, Turtle, Sea Lion, Shark) into a text encoder to generate textual embeddings $T_1, T_2, T_3, \dots, T_n$. An input image is processed by the image encoder to produce an image embedding I . Similarity scores between I and each $T_1, T_2, T_3, \dots, T_n$ are calculated, with the highest score (e.g., 0.48 for Turtle) indicating the model’s prediction for the image’s content.

these objects is challenging due to their varying sizes and their tendency to blend into the marine background, making this dataset particularly difficult for underwater object detection tasks.

- 3) **Underwater Object Counting Dataset:** This experiment uses the following dataset:

- a) **IOCfish5k [111]:** This dataset includes 5,637 annotated images with an average resolution of 1080×1920 pixels. It contains a total of 659,024 object counts, with minimum, average, and maximum object counts per image being 0, 117, and 2,371, respectively.

XII. SOTA METHODS FOR COMPARISON

We compared our AquaticCLIP model with a range of state-Of-The-Art (SOTA) methods across various underwater image analysis tasks.

- 1) **For the zero-shot classification task,** we evaluated the performance of AquaticCLIP against existing vision-language models (VLMs), including Frozen CLIP [97], Finetune CLIP [97], and prompt-based VLMs such as CoOp [143] and MAPLE [65], as well as GPT4V [129], BLIP2 [71], and MarineGPT [140]. CoOp and MAPLE were fine-tuned using the original source codes provided by their authors.
- 2) **For the supervised classification task,** we compared AquaticCLIP with models like ResNet-34/50/101 [51], ViT-S/B/L [17], BeiT [19], ConvNeXt [82], ConvNeXt [82] + Focal Loss (FL) [77], and ConvNeXt [82] + Class-Balanced (CB) [31]. To ensure a fair comparison, we followed the settings used in FishNet [64].
- 3) In the **underwater salient object segmentation task**, AquaticCLIP was evaluated against SOTA methods such as DUAL-SAM [135], SVAM-Net [58], CTDNet [137],

CDINet [132], SGL-KRN [124], TC-USOD [53], and MarineInst [149].

- 4) **For underwater instance segmentation,** we compared it with methods like Point Rend [67], SOLOv2 [122], QueryInst [41], Mask Transformer [63], Mask2Former [26], WaterMask R-CNN [111], and Cascade WaterMask R-CNN [111].
- 5) In the **underwater semantic segmentation task**, AquaticCLIP was compared to models including PSPNet [110], DeepLabv3 [23], SUIM-Net [56], and Mask2Former [26].
- 6) **For object detection and classification tasks,** we benchmarked AquaticCLIP against FasterRCNN [99], YOLOF [24], TOOD [42], MarineInst [149], MarineDet [48], and MRegionCLIP.
- 7) **For object counting tasks,** we compared it with existing SOTA methods, including KDMG [120], MPS [131], CLTR [76], and IOCFormer [111]. All methods were implemented using the official source codes provided by the original authors and evaluated on the same datasets for consistency.

XIII. SOTA METHODS TRAINING DETAILS

- **Supervised Classification Methods:** For a fair comparison, we used both ResNet-based and ViT-based architectures, including ResNets (ResNet-34, ResNet-50, ResNet-101) [51], ViT (ViT-S, ViT-B, ViT-L) [17], BeiT [19], and ConvNeXt [82]. The classification head was replaced with two fully connected (FC) layers, with a dropout rate of 0.5 applied to the first FC layer. All backbone networks were pre-trained on ImageNet and fine-tuned using the Adam optimizer over 100 epochs on each classification dataset’s training splits. The initial learning rates were set to 3×10^{-5} for the backbone networks and 3×10^{-4} for the classification heads, with the learning

rate halved every 20 epochs. To address the long-tail classification issue, we also evaluated focal loss [77] and class-balanced training [31] techniques and compared classification performance. For FishNet, all results are reported for fish family classification.

- **Supervised Object Detection and Classification Methods:** We evaluated five SOTA methods for object detection and classification: FastRCNN [99], YOLOF [24], and TOOD [42]. The backbone networks for these methods were pre-trained on the MS-COCO dataset [78], and we fine-tuned the complete models using the stochastic gradient descent (SGD) optimizer with a learning rate of 2.5×10^{-4} , momentum of 0.9, and weight decay of 1×10^{-4} for 80 epochs. Performance was reported on the official training and testing splits of the evaluated datasets to ensure fair comparison.

XIV. EVALUATION METRICS

We used various evaluation metrics suited to the specific aquatic imagery tasks.

- For **zero-shot classification tasks**, we reported accuracy and F_1 scores.
- For **zero-shot cross-modal retrieval tasks**, we evaluated performance using R@1, R@50, and R@200.
- In the **underwater salient object segmentation task**, metrics included S-measure (S_m), max E-measure (E_ϵ^{max}), max F-measure (max), and Mean Absolute Error (MAE) [53].
- For **instance segmentation** on the UIIS dataset, we used mAP, AP₅₀, AP₇₅, AP_S, AP_M, AP_L, AP_f, AP_h, and AP_r metrics [75].
- **Semantic segmentation** on the SUIM dataset was evaluated with IoU and F-measure scores.
- **Underwater object detection and classification tasks** are evaluated using mAP₅₀ metric.
- For **underwater object counting tasks**, we assessed performance using MAE and MSE metrics.

XV. ZERO-SHOT INFERENCE

Fig. 8 illustrates the zero-shot inference process, where normalized image and text embeddings are estimated and used for cosine similarity to assign class labels.

XVI. ZERO-SHOT CROSS-MODAL RETRIEVAL RESULTS

We evaluated zero-shot text-to-image and image-to-text retrieval tasks by finding the closest matches for each modality. Metrics such as R@1, R@50, and R@200 were used to check if the correct ground-truth pair was among the top matches. Table XII presents zero-shot cross-modal retrieval results on a dataset of 2K image-text pairs (excluded from the main 2M dataset), along with comparisons to four SOTA VLMs. AquaticCLIP outperformed all other methods by a significant margin, showcasing a strong alignment of cross-modal features across diverse visual and textual domains. MarineGPT ranked as the second-best performer.

TABLE VIII: Comparison of supervised object detection and classification results (mAP₅₀) across four datasets: FishNet, DeepFish, Brackish, and URPC. AquaticDet achieves the highest mAP₅₀ scores across all datasets, outperforming other SOTA methods. Specifically, AquaticDet scores 0.903 on FishNet, 0.891 on DeepFish, 0.877 on Brackish, and 0.837 on URPC. For FishNet, the performance is reported for common family classes.

Methods	FishNet	DeepFish	Brackish	URPC
FasterRCNN [99]	0.284	0.814	0.788	0.475
YOLOF [24]	0.672	0.806	0.813	0.511
TOOD [42]	0.811	0.766	0.805	0.507
MRegionCLIP	<u>0.867</u>	0.855	<u>0.842</u>	0.758
MarineInst [149]	<u>0.868</u>	0.854	0.841	<u>0.779</u>
MarineDet [48]	-	-	-	0.706
AquaticDet	0.903	0.891	0.877	0.837

TABLE IX: Supervised salient object segmentation results on the USOD10K dataset [53]. AquaticSAM outperforms other methods, including SVAM-Net [58], CTDNet [137], CDINet [132], SGL-KRN [124], TC-USOD [53], and MarineInst [149], achieving the best scores for S_m (0.943), E_ϵ^{max} (0.978), max F (0.942), and the lowest MAE (0.012), despite not using additional cues such as depth or edge information.

Methods	$S_m \uparrow$	$E_\epsilon^{max} \uparrow$	max F \uparrow	MAE \downarrow
SVAM-Net [58]	0.746	0.764	0.645	0.091
CTDNet [137]	0.908	0.953	0.907	0.028
CDINet [132]	0.704	0.864	0.736	0.090
SGL-KRN [124]	0.921	0.963	0.924	0.023
TC-USOD [53]	0.921	<u>0.968</u>	0.923	0.020
DUAL-SAM [135]	<u>0.923</u>	<u>0.968</u>	<u>0.931</u>	<u>0.018</u>
MarineInst [49]	0.910	0.941	0.887	0.025
AquaticSAM	0.943	0.978	0.942	0.012

XVII. UNDERWATER OBJECT DETECTION AND CLASSIFICATION RESULTS

For object detection and classification, we replaced ResNet-50 in MRegionCLIP with our pre-trained image encoder Ψ_v , naming the model AquaticDet. We applied the same fine-tuning settings as previously described. Table VIII shows object detection and classification results across four datasets, with AquaticDet achieving the highest mAP₅₀ scores across all datasets. This superior performance is attributed to pre-training on the 2M image-text paired dataset, which enabled AquaticDet to extract highly efficient visual features. The addition of the prompt-guided vision encoder and vision-guided text encoder, along with comprehensive captions at both image and instance levels, contributed to significant performance improvements even in challenging conditions.

XVIII. UNDERWATER SCENE SEGMENTATION RESULTS

For the supervised segmentation task, we used the SAM foundational model [66], which consists of a prompt encoder, an image encoder, and a lightweight mask decoder. In our fine-tuning experiments, we replaced SAM's original image encoder with the AquaticCLIP vision encoder Ψ_v , tuned explicitly for aquatic scenes. Other settings and implementation details followed SAM's original setup. We named our instance segmentation model AquaticSAM.

TABLE X: Comparison of instance segmentation results on the UIIS dataset [75] between AquaticSAM and other SOTA methods, all using ResNet101 as the backbone. AquaticSAM achieves the highest mAP (0.293), AP₅₀ (0.451), and AP_h (0.576), demonstrating superior performance in small object segmentation and overall accuracy.

Method	mAP \uparrow	AP ₅₀ \uparrow	AP ₇₅ \uparrow	AP _S \uparrow	AP _M \uparrow	AP _L \uparrow	AP _f \uparrow	AP _h \uparrow	AP _r \uparrow
Point Rend [67]	0.259	0.434	0.276	0.820	0.202	0.386	0.433	0.541	0.206
SOLov2 [122]	0.245	0.409	0.251	0.560	0.194	0.376	0.364	0.483	0.206
QueryInst [41]	0.260	0.428	0.273	0.820	0.217	0.351	0.433	0.541	0.206
Mask Transfiner [63]	0.246	0.421	0.260	0.720	0.194	0.361	0.438	0.263	0.198
Mask2Former [26]	0.257	0.380	0.277	0.630	0.189	0.381	0.411	0.519	0.231
WaterMask R-CNN [111]	<u>0.272</u>	<u>0.437</u>	0.293	0.900	<u>0.218</u>	<u>0.389</u>	0.463	0.548	0.209
Cascade WaterMask R-CNN [111]	0.271	0.429	<u>0.304</u>	0.830	0.210	<u>0.389</u>	<u>0.470</u>	<u>0.558</u>	<u>0.225</u>
MarineInst [149]	0.266	0.434	0.298	0.832	0.183	0.367	0.450	0.521	0.198
AquaticSAM	0.293	0.451	0.330	0.870	0.236	0.403	0.513	0.576	0.252

TABLE XI: Comparison of semantic segmentation results on the SUIM [56] dataset between AquaticSeg and other SOTA methods. AquaticSeg achieves the highest mIoU (0.881) and F-score (0.921), demonstrating superior performance in underwater semantic segmentation tasks.

Methods	mIoU \uparrow	F-score \uparrow
PSPNet [110]	0.774	0.760
DeepLabv3 [23]	0.791	0.812
SUIM-Net [56]	0.841	0.869
Mask2Former [26]	<u>0.855</u>	<u>0.896</u>
MarinsInst [149]	0.851	0.882
AquaticSeg	0.881	0.921

TABLE XII: Zero-shot cross-modal retrieval results (text-to-image and image-to-text) on the 2K dataset, presented as (text-to-image % | image-to-text %). AquaticCLIP outperforms other models, achieving the highest scores across R@1, R@50, and R@200, with top results of 22.29% | 21.91% at R@1, 62.31% | 63.20% at R@50, and 71.10% | 72.34% at R@200.

Methods	Our Dataset		
	R@1	R@50	R@200
Frozen CLIP	0.09 0.06	3.82 3.72	10.31 11.12
Finetune CLIP	<u>13.98</u> 12.79	43.42 42.31	61.34 60.74
GPT4V	8.19 8.17	35.10 34.13	40.61 41.12
MarineGPT	20.12 21.74	60.78 61.71	69.58 70.64
AquaticCLIP	22.29 21.91	62.31 63.20	71.10 72.34

- 1) Salient Object Segmentation Results:** Table IX shows the results of salient object segmentation on the USOD10K [53] dataset, where AquaticSAM consistently outperformed six existing SOTA methods across all metrics without relying on additional information such as edge or depth maps. While MarineInst also performed well, it lagged behind AquaticSAM.
- 2) Underwater Instance Segmentation Results:** In Table X, results for underwater instance segmentation on the UIIS [75] dataset indicate that AquaticSAM achieved the highest mAP, AP₅₀, and AP₇₅ scores, though it ranked second-best for small object segmentation (APS), with WaterMask R-CNN, its cascaded variant, and MarineInst also performing strongly.
- 3) Semantic Segmentation Results (Table XI):** For semantic segmentation, we introduced AquaticSeg, which incorporates a classification head fine-tuned for pixel-wise classification. Following the experimental protocols of MarineInst [149], these segmentation experiments used the training splits for each dataset. Table XI presents the semantic segmentation results on the SUIM [56] dataset, where AquaticSeg achieved the highest mIOU and F-score, with Mask2Former and MarineInst as the second- and third-best performers.

Overall, both AquaticSAM and AquaticSeg benefited from pre-training on the 2M aquatic images dataset. Additionally, the inclusion of both instance-level and image-level captions for language supervision significantly improved performance compared to using only image-level captions.

XIX. RESULTS OF OBJECT COUNTING IN UNDERWATER SCENES (TABLE XIII)

For object counting in underwater scenes, we employed a crowd localization transformer method [76], which includes a CNN-based backbone, a transformer encoder, a transformer decoder, and a nearest-neighbors matching component. In our experiments, we replaced the original backbone and encoder with the AquaticCLIP pre-trained vision encoder Ψ_v , specifically fine-tuned for aquatic environments. The remaining settings followed the original implementation. We refer to our object counting model as AquaticOC.

Table XIII presents object counting results on the IOCFish5K dataset [111], comparing AquaticOC with existing SOTA methods, including KDMG [120], MPS [131], CLTR [76], and IOCFormer [111]. All methods were implemented using the official source codes provided by the authors and evaluated on the same training and testing splits for fair comparison. AquaticOC achieved the best results in terms of MAE and MSE, showing a significant improvement over the baseline CLTR model, representing a strong advancement in efficient and accurate object counting for aquatic scenes.

XX. SUPERVISED AQUATICCLIP MODEL: LINEAR PROBE EVALUATIONS

We conducted linear probe experiments to evaluate the quality of visual features extracted by our AquaticCLIP model. Linear probing involves training a simple linear classifier, such as logistic regression, on top of features extracted from a pre-trained network, without fine-tuning the network's original weights. In this experiment, we applied logistic regression on

TABLE XIII: Object counting performance comparison of various SOTA methods on the IOCFish5k test set [111]. AquaticOC achieves the best results, with the lowest Mean Absolute Error (MAE: 13.50) and Mean Squared Error (MSE: 36.10), outperforming other methods.

Method	MAE ↓	MSE ↓
KDMG [120]	0.227	0.499
MPS [131]	0.335	0.550
CLTR [76]	0.180	0.419
IOCFormer [111]	0.171	0.412
AquaticOC	0.135	0.361

pre-extracted features from the AquaticCLIP vision encoder to assess the discriminative power and quality of the representations.

For downstream classification tasks across various marine vision datasets, we used logistic regression, following practices recommended by the self-supervised learning community [22], [44]. We set the ℓ_2 regularization coefficient λ to $\frac{100}{MC}$, where M is the embedding dimension and C is the number of classes, and used the L-BFGS solver with a maximum of 1,000 iterations [145]. All downstream classification tasks were evaluated using stratified train-test splits or official splits when available.

XXI. AQUATICVISION: VISION-ONLY MODEL PRE-TRAINING DETAILS

For comparison, we also developed a vision-only model named AquaticVision, pre-trained using a self-supervised learning approach. For large-scale visual pre-training on our 2M aquatic imagery dataset, we employed DINOv2 [92], a SOTA self-supervised learning method based on student-teacher knowledge distillation for pre-training large ViT architectures.

We used a ViT-B architecture with 12 layers, 12 heads, an 8×8 patch size, a feed-forward network layer as MLP, GELU activation, an embedding dimension of 768, and a dropout rate of 0.1 [38]. The model was trained with the AdamW optimizer with hyperparameter β set to (0.9, 0.999), a batch size of 3072, and a maximum of 125,000 iterations.

XXII. VISUAL RESULTS

Figures 9-12 display sample results from various downstream analysis tasks.

XXIII. WHY AQUATICCLIP PERFORMANCE IS BETTER?

While MarineInst [149] effectively performs instance segmentation and captioning in marine environments, AquaticCLIP addresses some of its limitations, primarily focusing on enhancing accuracy and addressing semantic understanding in a zero-shot setting.

1. Addressing Over-Segmentation and Partial-Segmentation: MarineInst acknowledges challenges with over-segmentation and partial-segmentation in complex marine images. AquaticCLIP, while not directly addressing MarineInst, tackles similar issues through its prompt-guided

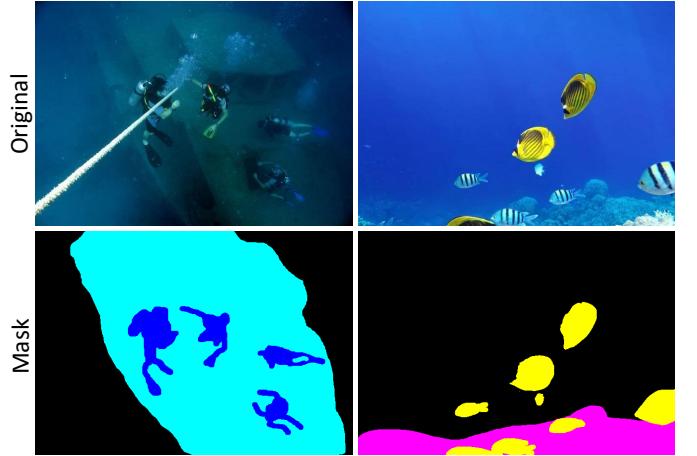


Fig. 9: Semantic segmentation results on the SUIM dataset [56]. The top row shows original underwater images, while the bottom row displays corresponding segmentation masks generated by AquaticSeg, highlighting detected objects such as divers and fish with distinct color-coded regions.



Fig. 10: Instance segmentation results on the UIIS dataset [75]. The top row shows original underwater scenes, while the bottom row presents segmented images generated by AquaticSAM, highlighting individual objects like fish, corals, and sharks with distinct boundaries for each detected instance.



Fig. 11: Underwater object detection results on the FishNet dataset [64]. The images show detected objects, including a ray, crocodile, and various fish, each marked with bounding boxes and confidence scores, demonstrating AquaticDet's capability to accurately identify and label underwater species.

vision encoder (PGVE) and a dedicated textual description cleaning module (TDCM).

- The PGVE uses learned visual prompts to aggregate patch features, prioritizing semantically similar regions and

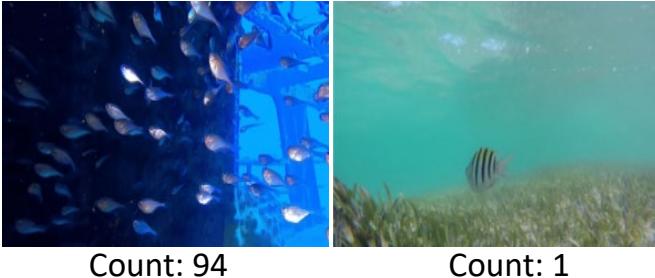


Fig. 12: Underwater object counting results on the IOCfish5k dataset [111]. The images show object counts generated by AquaticOC, with the left image counting 94 fish and the right image counting 1 fish, demonstrating the model’s ability to accurately detect and count objects in various underwater scenes.

creating a more meaningful image-level representation. This focus on relevant patches helps reduce the likelihood of segmenting extraneous regions, thus addressing over-segmentation.

- The TDCM refines the automatically generated textual descriptions by identifying and retaining the most semantically relevant keywords. This cleaning process ensures higher-quality textual descriptions, which in turn, contributes to more accurate instance identification and segmentation, potentially reducing instances of partial segmentation.

2. Enhancing Semantic Understanding in a Zero-Shot Setting:

AquaticCLIP excels in zero-shot classification tasks, achieving high accuracy in recognizing marine species and coral categories not encountered during training. This is a significant advantage over MarineInst, which relies heavily on its pre-defined training categories for semantic understanding.

- AquaticCLIP leverages the knowledge from pre-trained language models like MarineGPT, GPT4V, and BLIP2 to generate textual descriptions for the images. This allows for a broader and more nuanced understanding of the visual content, even for novel classes.
- The use of a vision-guided text encoder (VGTE) further enhances semantic understanding by incorporating visual context from the image into the text features. This alignment between visual and textual modalities results in richer semantic representations, crucial for accurate zero-shot classification.

In essence, while both models contribute significantly to marine image analysis, AquaticCLIP’s focus on precise feature aggregation through PGVE, refined textual descriptions with TDCM, and enhanced semantic understanding with VGTE in a zero-shot setting addresses some limitations inherent in MarineInst’s approach.

3. AquaticCLIP’s Efficient Training and Comparison to MarineInst:

- AquaticCLIP is trained on a dataset of 2 million image-text pairs. This dataset focuses specifically on aquatic environments and is curated from diverse sources like YouTube, marine biology texts, online repositories, and social media etc. The key is that while AquaticCLIP uses real aquatic images, the textual descriptions are pseudo-

generated, meaning they are created automatically rather than manually annotated. This process involves using pre-trained language models like MarineGPT and applying a cleaning module to refine the generated text.

- MarineInst, on the other hand, is trained on a dataset called MarineInst20M, containing 2.4 million marine images with 20M object instances. This dataset includes images from public underwater datasets, manually collected images, and public internet images. Unlike AquaticCLIP, MarineInst20M utilizes a mixture of human-annotated and automatically generated instance masks.

Although trained on a 2M image-text paired dataset with both groundtruth and additional textual descriptions, AquaticCLIP demonstrates comparable and often superior performance to MarineInst across various marine vision tasks. This efficiency suggests that AquaticCLIP’s architecture, particularly its use of prompt-guided vision encoding and vision-guided text encoding, effectively leverages the available data for robust and accurate analysis.

REFERENCES

- [1] “,” Available on Reelfile survey.
- [2] “,” Available on Reeflex.
- [3] “Aquarium dataset,” Available on Aquarium.
- [4] “Flicker Images.”
- [5] “Getty images.”
- [6] “Hk reef fish Images.”
- [7] “Marine Animal Images,” Available on Kaggle.
- [8] “Sea Animals Image Dataset,” Available on Kaggle.
- [9] “Shutterstock Image.”
- [10] “Underwater trash detection dataset,” Available on Roboflow.
- [11] Available on EOL, 2018.
- [12] “Ozfish dataset - machine learning dataset for baited remote underwater video stations,” 2020.
- [13] “URPC dataset,” Available on Kaggle, 2020.
- [14] “Oceanic life dataset,” Available on Kaggle, 2023.
- [15] B. Ai, X. Liu, Z. Wen, L. Wang, H. Ma, and G. Lv, “A novel coral reef classification method combining radiative transfer model with deep learning,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2024.
- [16] B. Alawode, Y. Guo, M. Ummar, N. Werghi, J. Dias, A. Mian, and S. Javed, “Utb180: A high-quality benchmark for underwater tracking,” in *Proceedings of the Asian Conference on Computer Vision*, 2022, pp. 3326–3342.
- [17] D. Alexey, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv: 2010.11929*, 2020.
- [18] K. M. Babu, D. Bentall, D. T. Ashton, M. Pukowski, W. Fantham, H. T. Lin, N. P. Tuckey, M. Wellenreuther, and L. K. Jesson, “Computer vision in aquaculture: a case study of juvenile fish counting,” *Journal of the Royal Society of New Zealand*, vol. 53, no. 1, pp. 52–68, 2023.
- [19] H. Bao, L. Dong, S. Piao, and F. Wei, “Beit: Bert pre-training of image transformers,” *arXiv preprint arXiv:2106.08254*, 2021.
- [20] A. B. Burguera, F. Bonin-Font, D. Chatzivangelou, M. V. Fernandez, and J. Aguzzi, “Deep learning for detection and counting of nephrops norvegicus from underwater videos,” *ICES Journal of Marine Science*, p. fsae089, 2024.
- [21] L. Cai, N. E. McGuire, R. Hanlon, T. A. Mooney, and Y. Girdhar, “Semi-supervised visual tracking of marine animals using autonomous underwater vehicles,” *International Journal of Computer Vision*, vol. 131, no. 6, pp. 1406–1427, 2023.
- [22] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, “Emerging properties in self-supervised vision transformers,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 9650–9660.
- [23] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 834–848, 2018.

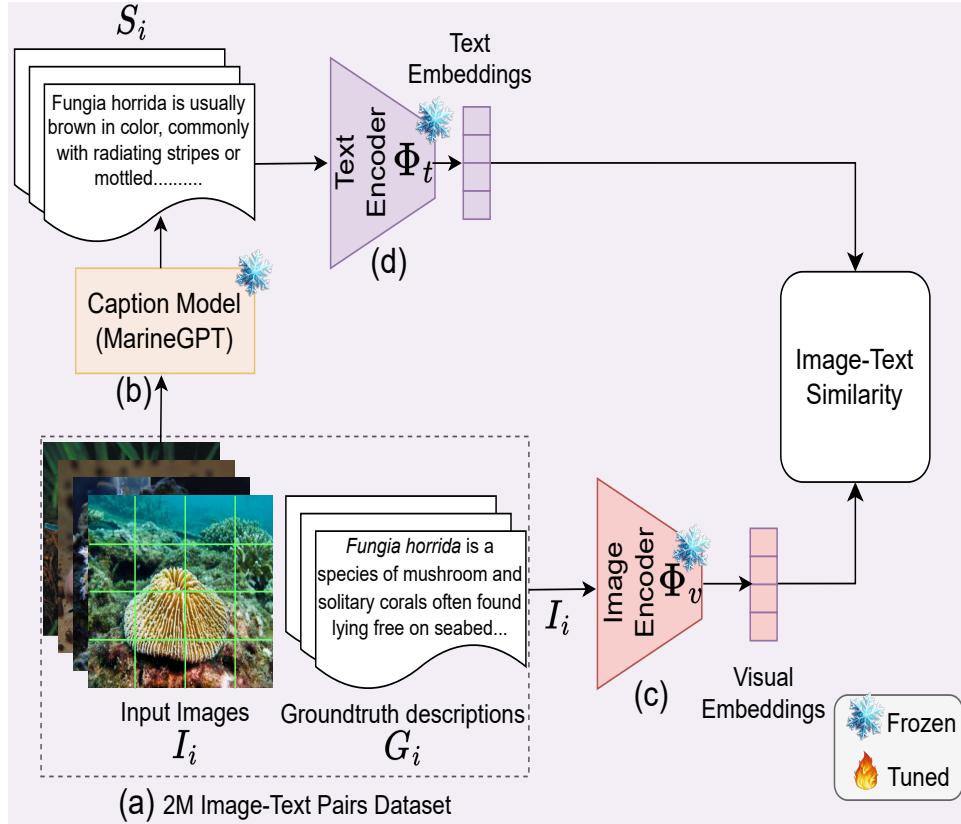


Fig. 13: Schematic illustration of the Frozen CLIP. Please refer to Table 1 in the main manuscript.

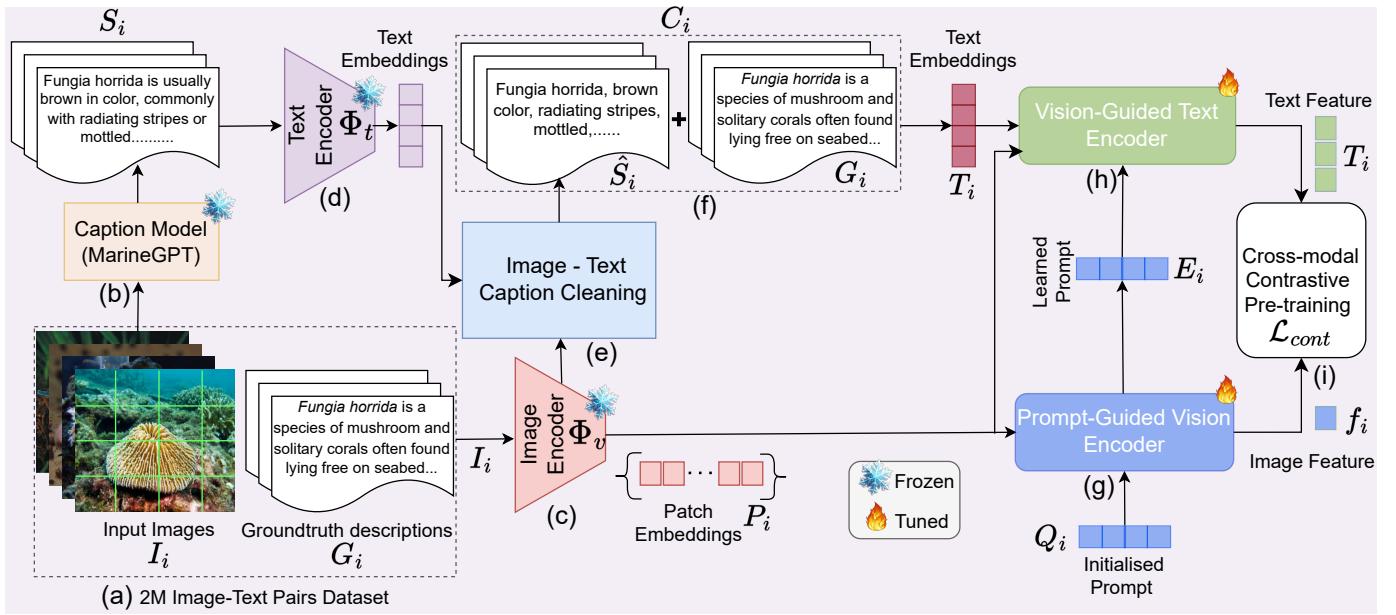


Fig. 14: Schematic illustration of the AquaticCLIP₁. Please refer to Table 1 in the main manuscript.

- [24] Q. Chen, Y. Wang, T. Yang, X. Zhang, J. Cheng, and J. Sun, “You only look one-level feature,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 13 039–13 048.
- [25] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, “A simple framework for contrastive learning of visual representations,” in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [26] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, “Masked-attention mask transformer for universal image segmentation,” in *CVPR*, 2022.
- [27] Y. Cheng, J. Zhu, M. Jiang, J. Fu, C. Pang, P. Wang, K. Sankaran, O. Onabola, Y. Liu, D. Liu, and Y. Bengio, “Flow: A dataset and benchmark for floating waste detection in inland waters,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 10 953–10 962.
- [28] C. Chin, “Marine fouling images.” 2019.
- [29] Y. Cho, H. Jang, R. Malav, G. Pandey, and A. Kim, “Underwater image dehazing via unpaired image-to-image translation,” *International Journal of Control, Automation and Systems*, vol. 18, pp. 605–614,

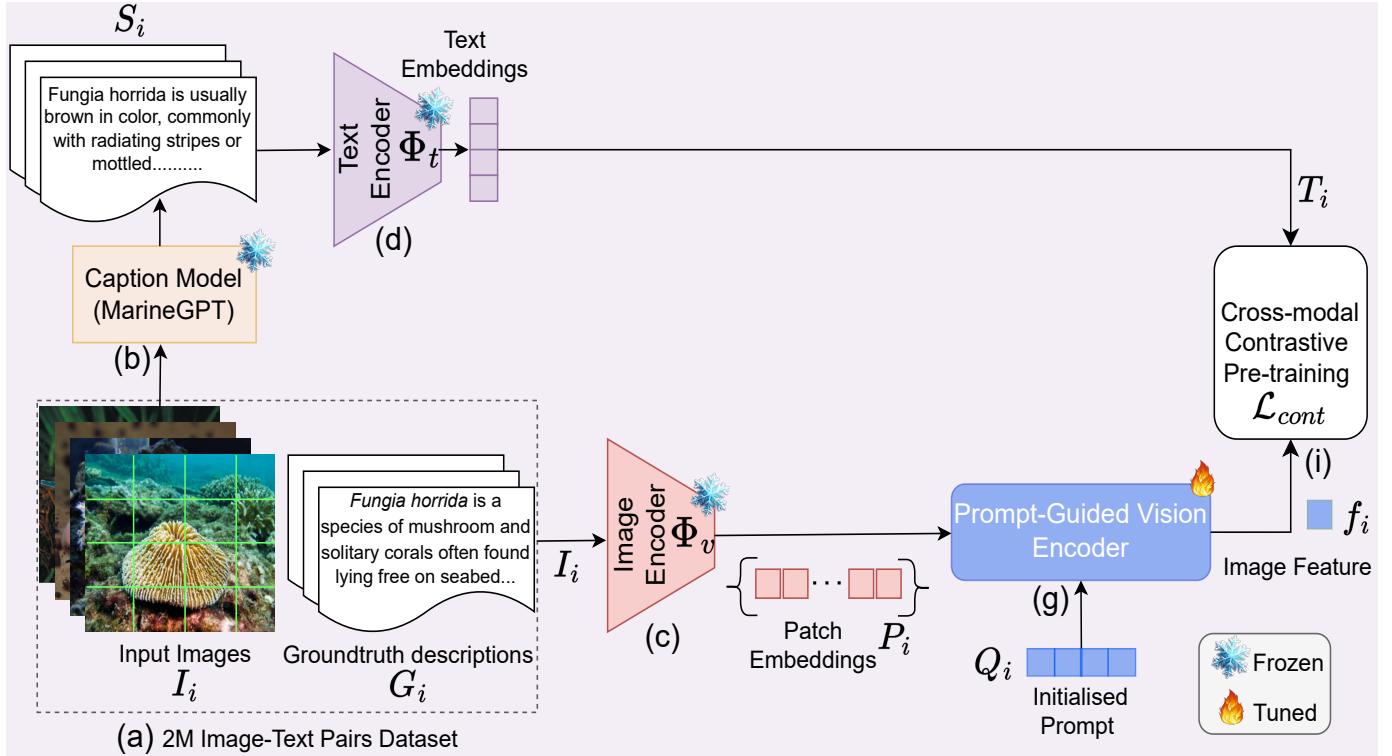


Fig. 15: Schematic illustration of the AquaticCLIP₂. Please refer to Table 1 in the main manuscript.

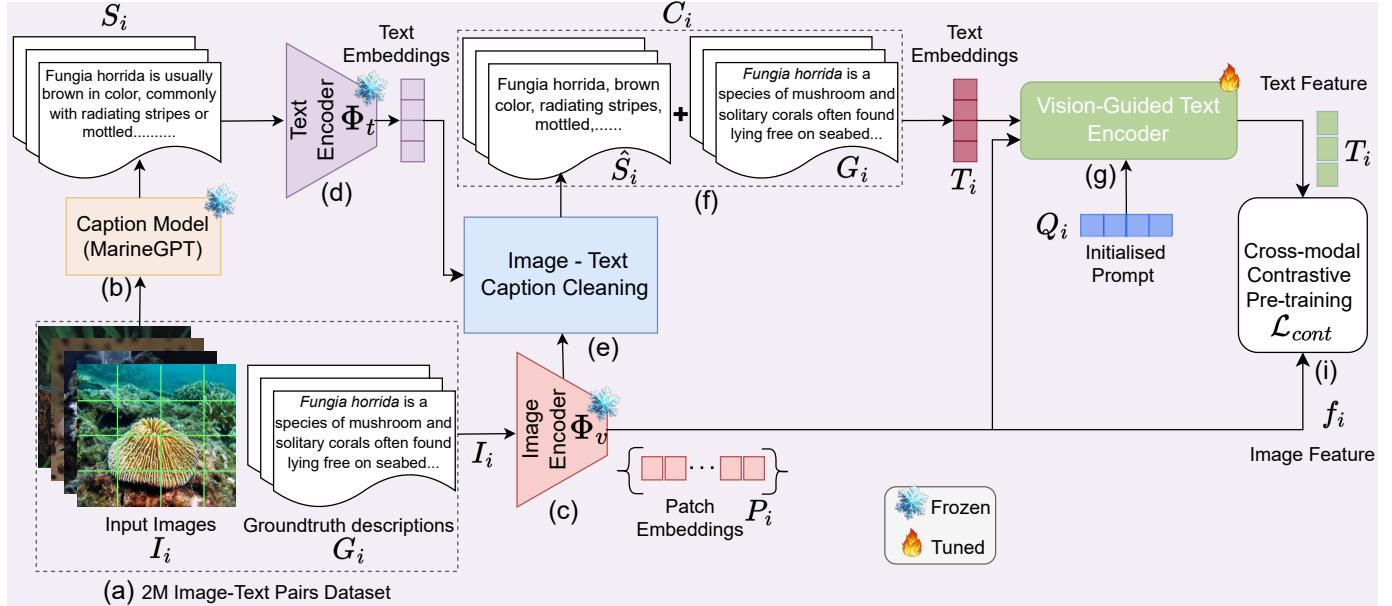


Fig. 16: Schematic illustration of the AquaticCLIP₃. Please refer to Table 1 in the main manuscript.

2020.

- [30] C. Clark and S. Divvala, "Pdffigures 2.0: Mining figures from research papers," in *Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries*, 2016, pp. 143–152.
- [31] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 9268–9277.
- [32] K. Dell, "Ocean acidification's impact on marine ecosystems," *National Geographic*, May 2024, accessed: November 9, 2024.
- [33] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPR Workshops)*. Los Alamitos, CA, USA: IEEE Computer Society, jun 2009, pp. 248–255.
- [34] C. Desai, B. S. S. Reddy, R. A. Tabib, U. Patil, and U. Mudenagudi, "Aquagan: Restoration of underwater images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 296–304.
- [35] A. Dighe, "Corals classification," Available on Kaggle.
- [36] J. Doe, "The wonders of coral reefs," *National Geographic*, vol. 245, no. 4, pp. 56–71, April 2023.
- [37] S. C. Doney, M. Ruckelshaus, J. Emmett Duffy, J. P. Barry, F. Chan, C. A. English, H. M. Galindo, J. M. Grebmeier, A. B. Hollowed,

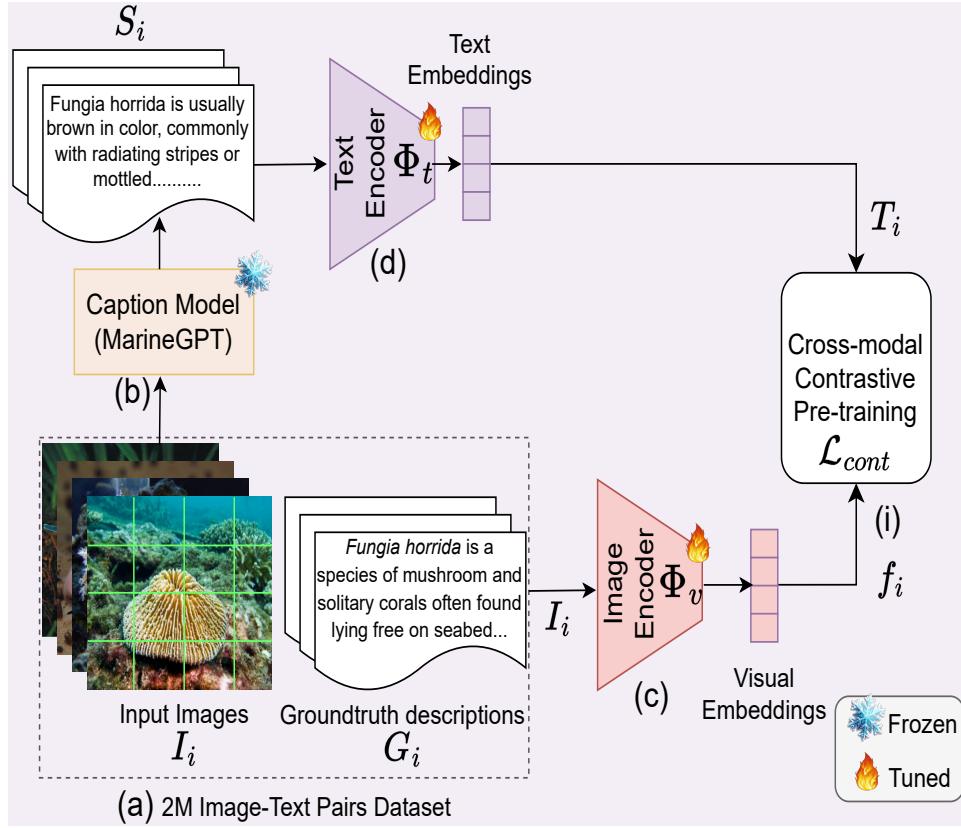


Fig. 17: Schematic illustration of the Finetuned CLIP. Please refer to Table 1 in the main manuscript.

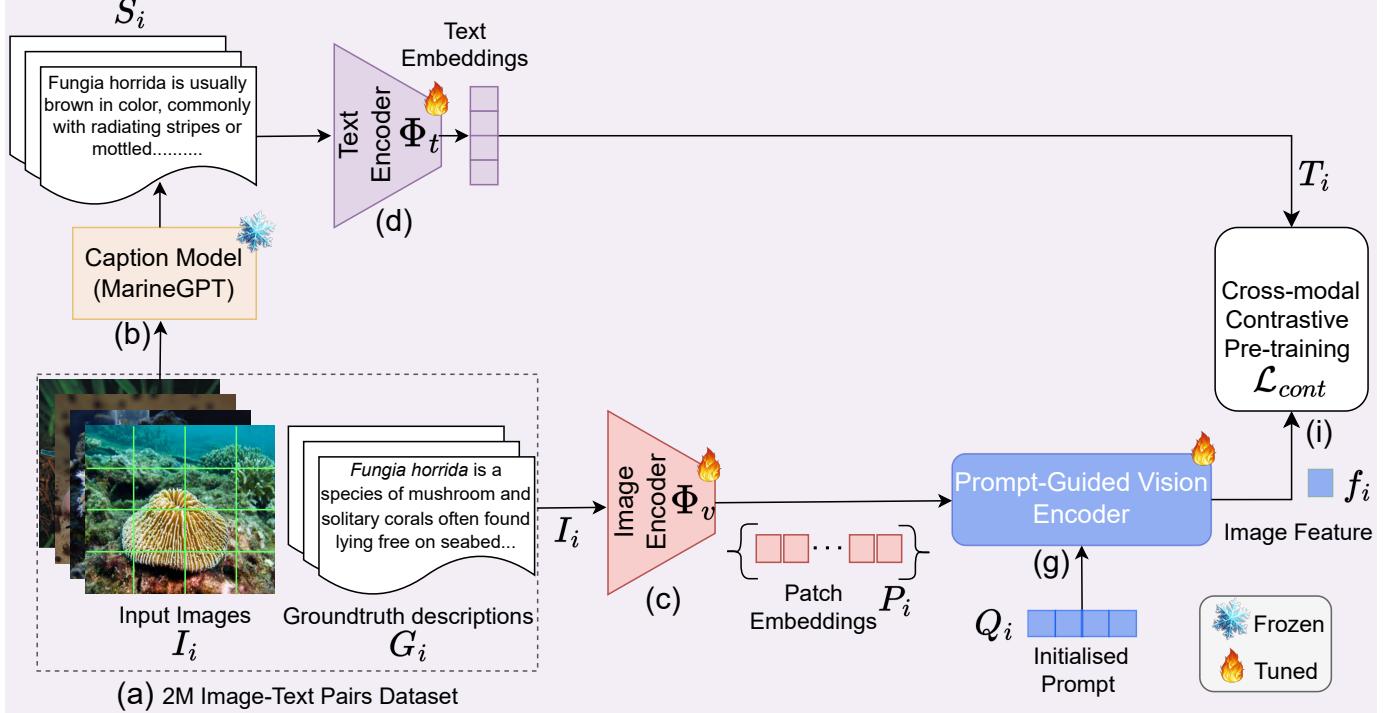


Fig. 18: Schematic illustration of the AquaticCLIP4. Please refer to Table 1 in the main manuscript.

- N. Knowlton *et al.*, “Climate change impacts on marine ecosystems,” *Annual review of marine science*, vol. 4, no. 1, pp. 11–37, 2012.
- [38] A. Dosovitskiy, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [39] C. M. Duarte, S. Agusti, E. Barbier, G. L. Britten, J. C. Castilla, J.-P. Gattuso, R. W. Fulweiler, T. P. Hughes, N. Knowlton, C. E. Lovelock *et al.*, “Rebuilding marine life,” *Nature*, vol. 580, no. 7801, pp. 39–51, 2020.
- [40] J. Evers, “The decline of australia’s great barrier reef,” *National Geographic*, September 2023, accessed: November 9, 2024.

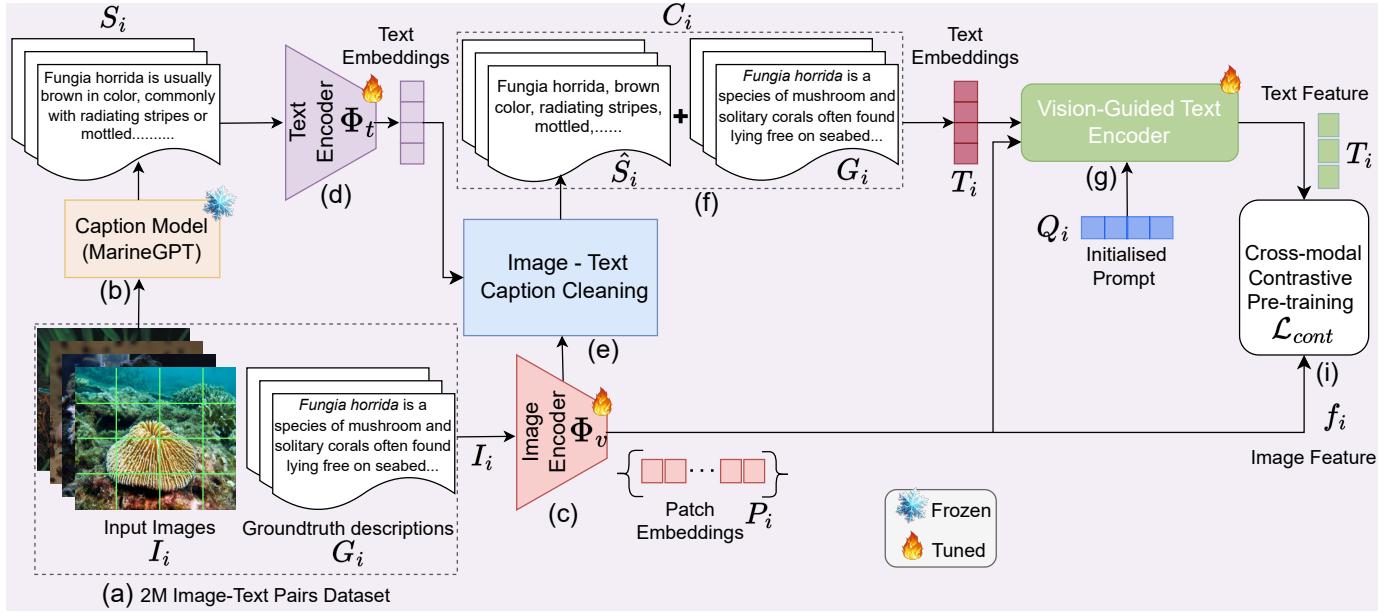


Fig. 19: Schematic illustration of the AquaticCLIP₅. Please refer to Table 1 in the main manuscript.

- [41] Y. Fang, S. Yang, X. Wang, Y. Li, C. Fang, Y. Shan, B. Feng, and W. Liu, “Instances as queries,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 6910–6919.
- [42] C. Feng, Y. Zhong, Y. Gao, M. R. Scott, and W. Huang, “Tood: Task-aligned one-stage object detection,” in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE Computer Society, 2021, pp. 3490–3499.
- [43] M. Fulton, J. Hong, M. J. Islam, and J. Sattar, “Robotic detection of marine litter using deep visual detection models,” in *2019 International Conference on Robotics and Automation (ICRA)*, 2019, pp. 5752–5758.
- [44] J. Geiping, Q. Garrido, P. Fernandez, A. Bar, H. Pirsiavash, Y. LeCun, and M. Goldblum, “A cookbook of self-supervised learning,” *arXiv preprint arXiv:2304.12210*, 2023.
- [45] S. P. González-Sabbagh and A. Robles-Kelly, “A survey on underwater computer vision,” *ACM Computing Surveys*, vol. 55, no. 13s, pp. 1–39, 2023.
- [46] K. Grorud-Colvert, J. Sullivan-Stack, C. Roberts, V. Constant, B. Horta e Costa, E. P. Pike, N. Kingston, D. Laffoley, E. Sala, J. Claudet *et al.*, “The mpa guide: A framework to achieve global goals for the ocean,” *Science*, vol. 373, no. 6560, p. eabf0861, 2021.
- [47] G. Ha, “Coral species classification dataset,” Available on Roboflow, 2022.
- [48] L. Haixin, Z. Ziqiang, M. Zeyu, and S.-K. Yeung, “Marinedet: Towards open-marine object detection,” *arXiv preprint arXiv:2310.01931*, 2023.
- [49] B. S. Halpern, S. Walbridge, K. A. Selkoe, C. V. Kappel, F. Micheli, C. d’Agrosa, J. F. Bruno, K. S. Casey, C. Ebert, H. E. Fox *et al.*, “A global map of human impact on marine ecosystems,” *science*, vol. 319, no. 5865, pp. 948–952, 2008.
- [50] Z. Hao, J. Qiu, H. Zhang, G. Ren, and C. Liu, “Umotma: Underwater multiple object tracking with memory aggregation,” *Frontiers in Marine Science*, vol. 9, p. 1071618, 2022.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [52] J. Hong, M. Fulton, and J. Sattar, “Trashcan: A semantically-segmented dataset towards visual detection of marine debris,” *CoRR*, vol. abs/2007.08097, 2020.
- [53] L. Hong, X. Wang, G. Zhang, and M. Zhao, “Usod10k: A new benchmark dataset for underwater salient object detection,” *IEEE Transactions on Image Processing*, pp. 1–1, 2023.
- [54] Y. Hu, K. Wang, X. Zhao, H. Wang, and Y. Li, “Underwater image restoration based on convolutional neural network,” in *Asian conference on machine learning*. PMLR, 2018, pp. 296–311.
- [55] S. Huang, K. Wang, H. Liu, J. Chen, and Y. Li, “Contrastive semi-supervised learning for underwater image restoration via reliable bank,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 18 145–18 155.
- [56] M. J. Islam, C. Edge, Y. Xiao, P. Luo, M. Mehtaz, C. Morse, S. S. Enan, and J. Sattar, “Semantic Segmentation of Underwater Imagery: Dataset and Benchmark,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE/RSJ, 2020.
- [57] M. J. Islam, Y. Xia, and J. Sattar, “Fast underwater image enhancement for improved visual perception,” *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3227–3234, 2020.
- [58] M. Jahidul Islam, R. Wang, and J. Sattar, “Svam: Saliency-guided visual attention modeling by autonomous underwater robots,” *arXiv e-prints*, pp. arXiv–2011, 2020.
- [59] A. Jalal, A. Salman, A. Mian, M. Shortis, and F. Shafait, “Fish detection and species classification in underwater environments using deep learning with temporal information,” *Ecological Informatics*, vol. 57, p. 101088, 2020.
- [60] S. Jennings and M. J. Kaiser, “The effects of fishing on marine ecosystems,” in *Advances in marine biology*. Elsevier, 1998, vol. 34, pp. 201–352.
- [61] K. Katija, E. Orenstein, B. Schlining, L. Lundsten, K. Barnard, G. Sainz, O. Boulais, M. Cromwell, E. Butler, B. Woodward, and K. L. Bell, “FathomNet: A global image database for enabling artificial intelligence in the ocean,” *Sci. Rep.*, vol. 12, no. 1, pp. 1–14, 2022.
- [62] J. Kay and M. Merrifield, “The fishnet open images database: A dataset for fish detection and fine-grained categorization in fisheries,” *arXiv preprint arXiv:2106.09178*, 2021.
- [63] L. Ke, M. Danelljan, X. Li, Y.-W. Tai, C.-K. Tang, and F. Yu, “Mask transfiner for high-quality instance segmentation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 4412–4421.
- [64] F. F. Khan, X. Li, A. J. Temple, and M. Elhoseiny, “Fishnet: A large-scale dataset and benchmark for fish recognition, detection, and functional trait prediction,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 20 496–20 506.
- [65] M. U. Khattak, H. Rasheed, M. Maaz, S. Khan, and F. S. Khan, “Maple: Multi-modal prompt learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 19 113–19 122.
- [66] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollar, and R. Girshick, “Segment anything,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 4015–4026.
- [67] A. Kirillov, Y. Wu, K. He, and R. Girshick, “Pointrend: Image segmentation as rendering,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [68] S.-H. Lee and M.-H. Oh, “Detection and tracking of underwater fish using the fair multi-object tracking model: A comparative analysis

- of yolov5s and dla-34 backbone models,” *Applied Sciences*, vol. 14, no. 16, p. 6888, 2024.
- [69] J. Li, K. A. Skinner, R. M. Eustice, and M. Johnson-Roberson, “Watrgan: Unsupervised generative network to enable real-time color correction of monocular underwater images,” *IEEE Robotics and Automation letters*, vol. 3, no. 1, pp. 387–394, 2017.
- [70] J. Li, W. Xu, L. Deng, Y. Xiao, Z. Han, and H. Zheng, “Deep learning for visual recognition and detection of aquatic animals: A review,” *Reviews in Aquaculture*, vol. 15, no. 2, pp. 409–433, 2023.
- [71] J. Li, D. Li, S. Savarese, and S. Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” in *International conference on machine learning*. PMLR, 2023, pp. 19 730–19 742.
- [72] X. Li, Y. Huang, Z. He, Y. Wang, H. Lu, and M.-H. Yang, “Citetracker: Correlating image and text for visual tracking,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 9974–9983.
- [73] Y. Li, F. Liang, L. Zhao, Y. Cui, W. Ouyang, J. Shao, F. Yu, and J. Yan, “Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm,” *arXiv preprint arXiv:2110.05208*, 2021.
- [74] Y. Li, B. Wang, Y. Li, Z. Liu, W. Huo, Y. Li, and J. Cao, “Underwater object tracker: Uostrack for marine organism grasping of underwater vehicles,” *Ocean Engineering*, vol. 285, p. 115449, 2023.
- [75] S. Lian, H. Li, R. Cong, S. Li, W. Zhang, and S. Kwong, “Watermask: Instance segmentation for underwater imagery,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2023, pp. 1305–1315.
- [76] D. Liang, W. Xu, and X. Bai, “An end-to-end transformer model for crowd localization,” *European Conference on Computer Vision*, 2022.
- [77] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [78] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.
- [79] C. Liu, Z. Wang, S. Wang, T. Tang, Y. Tao, C. Yang, H. Li, X. Liu, and X. Fan, “A new dataset, poisson gan and aquanet for underwater object grabbing,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 5, pp. 2831–2844, 2022.
- [80] J. Liu, A. T. Becerra, J. F. Bienvenido-Barcena, X. Yang, Z. Zhao, and C. Zhou, “Cffi-vit: Enhanced vision transformer for the accurate classification of fish feeding intensity in aquaculture,” *Journal of Marine Science and Engineering*, vol. 12, no. 7, p. 1132, 2024.
- [81] Y. Liu, D. An, Y. Ren, J. Zhao, C. Zhang, J. Cheng, J. Liu, and Y. Wei, “Dp-fishnet: Dual-path pyramid vision transformer-based underwater fish detection network,” *Expert Systems with Applications*, vol. 238, p. 122018, 2024.
- [82] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 11 976–11 986.
- [83] I. Loshchilov, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [84] J. Lu, N. Li, S. Zhang, Z. Yu, H. Zheng, and B. Zheng, “Multi-scale adversarial networks for underwater image restoration,” *Optics & Laser Technology*, vol. 110, pp. 105–113, 2019.
- [85] M. Manipambil, C. Vorster, D. Molloy, N. Murphy, K. McGuinness, and N. E. O’Connor, “Enhancing clip with gpt-4: Harnessing visual descriptions as prompts,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 262–271.
- [86] J. R. Merrick, “Australasian freshwater fish faunas: diversity, interrelationships, radiations and conservation,” in *Evolution and biogeography of Australasian vertebrates*. Auscipub, 2006, pp. 195–224.
- [87] Z. Michel, “The rich biodiversity of ocean ecosystems,” *National Geographic*, December 2022, accessed: November 9, 2024.
- [88] M. Modasshir, S. Rahman, O. Youngquist, and I. Rekleitis, “Coral identification and counting with an autonomous underwater vehicle,” in *2018 IEEE international conference on robotics and biomimetics (ROBIO)*. IEEE, 2018, pp. 524–529.
- [89] L. Mohan, “Innovative approaches to coral conservation,” *National Geographic*, June 2023, accessed: November 9, 2024.
- [90] Y. Mu, Q. Zhang, M. Hu, W. Wang, M. Ding, J. Jin, B. Wang, J. Dai, Y. Qiao, and P. Luo, “Embodiedgpt: Vision-language pre-training via embodied chain of thought,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [91] M. F. Naeem, M. G. Z. A. Khan, Y. Xian, M. Z. Afzal, D. Stricker, L. Van Gool, and F. Tombari, “I2mvformer: Large language model generated multi-view document supervision for zero-shot image classification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 15 169–15 179.
- [92] M. Oquab, T. Dariset, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby *et al.*, “Dinov2: Learning robust visual features without supervision,” 2023.
- [93] K. Panetta, L. Kezebou, V. Oludare, and S. Agaian, “Comprehensive underwater object tracking benchmark dataset and underwater image enhancement with gan,” *IEEE Journal of Oceanic Engineering*, vol. 47, no. 1, pp. 59–75, 2022.
- [94] M. Pedersen, J. Bruslund Haurum, R. Gade, and T. B. Moeslund, “Detection of marine animals in a new underwater dataset with varying visibility,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- [95] L. Peng, C. Zhu, and L. Bian, “U-shape transformer for underwater image enhancement,” *IEEE Transactions on Image Processing*, vol. 32, pp. 3066–3079, 2023.
- [96] C. Prathima, C. Silpa, A. Charitha, G. Harshitha, C. Sai Charan, and G. R. Sailendra, “Detecting and recognizing marine animals using advanced deep learning models,” in *2024 International Conference on Expert Clouds and Applications (ICOECA)*, 2024, pp. 950–955.
- [97] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [98] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [99] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.
- [100] A. Saleh, I. H. Laradji, D. A. Konovalov, M. Bradley, D. Vazquez, and M. Sheaves, “A realistic fish-habitat dataset to evaluate algorithms for underwater visual analysis,” *Scientific Reports*, vol. 10, no. 1, p. 14671, 2020.
- [101] A. Saleh, M. Sheaves, D. Jerry, and M. R. Azghadi, “Transformer-based self-supervised fish segmentation in underwater videos,” *arXiv preprint arXiv:2206.05390*, 2022.
- [102] A. Saleh, M. Sheaves, and M. Rahimi Azghadi, “Computer vision and deep learning for fish classification in underwater habitats: A survey,” *Fish and Fisheries*, vol. 23, no. 4, pp. 977–999, 2022.
- [103] F. Sammani, T. Mukherjee, and N. Deligiannis, “NLx-gpt: A model for natural language explanations in vision and vision-language tasks,” in *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 8322–8332.
- [104] R. Schodde, *Zoological Catalogue of Australia: Aves (Columbidae to Coraciidae)*. CSIRO PUBLISHING, 1997, vol. 37.
- [105] X. Shao, H. Chen, K. Magson, J. Wang, J. Song, J. Chen, and J. Sasaki, “Deep learning for multi-label classification of coral conditions in the indo-pacific via underwater photogrammetry,” *arXiv preprint arXiv:2403.05930*, 2024.
- [106] J. J. Shelley, A. Delaval, and M. C. Le Feuvre, “A revision of the grunter genus syncomistes (teleostei, terapontidae, syncomistes) with descriptions of seven new species from the kimberley region, northwestern australia,” *Zootaxa*, vol. 4367, no. 1, pp. 1–103, 2017.
- [107] G. Si, Y. Xiao, B. Wei, L. B. Bullock, Y. Wang, and X. Wang, “Token-selective vision transformer for fine-grained image recognition of marine organisms,” *Frontiers in Marine Science*, vol. 10, p. 1174347, 2023.
- [108] S. A. Siddiqui, A. Salman, M. I. Malik, F. Shafait, A. Mian, M. R. Shortis, and E. S. Harvey, “Automatic fish species classification in underwater videos: exploiting pre-trained deep neural network models to compensate for limited labelled data,” *ICES Journal of Marine Science*, vol. 75, no. 1, pp. 374–389, 2018.
- [109] A. Singh, R. Hu, V. Goswami, G. Couairon, W. Galuba, M. Rohrbach, and D. Kiela, “Flava: A foundational language and vision alignment model,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 15 638–15 650.
- [110] Q. Song, C. Wang, Z. Jiang, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, and Y. Wu, “Rethinking counting and localization in crowds: A purely point-based framework,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 3365–3374.
- [111] G. Sun, Z. An, Y. Liu, C. Liu, C. Sakaridis, D.-P. Fan, and L. Van Gool, “Indiscernible object counting in underwater scenes,” in *Proceedings*

- of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 13 791–13 801.
- [112] Z. Sun, Q. Ke, H. Rahmani, M. Bennamoun, G. Wang, and J. Liu, “Human action recognition from various data modalities: A review,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 3, pp. 3200–3225, 2022.
- [113] S. Thornton, “Understanding coral bleaching and ocean health,” *National Geographic*, July 2024, accessed: November 9, 2024.
- [114] S. Thornton and L. J. Richardson, “The vital role of coral reefs and the threats they face,” *National Geographic*, July 2024, accessed: November 9, 2024.
- [115] Y. Tian, D. Krishnan, and P. Isola, “Contrastive multiview coding,” in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*. Springer, 2020, pp. 776–794.
- [116] O. Ulucan, D. Karakaya, and M. Turkan, “A large-scale dataset for fish segmentation and classification,” in *2020 Innovations in Intelligent Systems and Applications Conference (ASYU)*. IEEE, 2020, pp. 1–5.
- [117] M. Ummar, F. A. Dharejo, B. Alawode, T. Mahbub, M. J. Piran, and S. Javed, “Window-based transformer generative adversarial network for autonomous underwater image enhancement,” *Engineering Applications of Artificial Intelligence*, vol. 126, p. 107069, 2023.
- [118] R. Van Der Laan, W. N. Eschmeyer, and R. Fricke, “Family-group names of recent fishes,” *Zootaxa*, vol. 3882, no. 1, pp. 1–230, 2014.
- [119] J. Veron, M. Stafford-Smith, E. Turak, and L. DeVantier, “Corals of the world. accessed 12 mar 2023, version 0.01,” 2016.
- [120] J. Wan, Q. Wang, and A. B. Chan, “Kernel-based density map generation for dense object counting,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 3, pp. 1357–1370, 2022.
- [121] N. Wang, T. Chen, S. Liu, R. Wang, H. R. Karimi, and Y. Lin, “Deep learning-based visual detection of marine organisms: A survey,” *Neurocomputing*, vol. 532, pp. 1–32, 2023.
- [122] X. Wang, R. Zhang, T. Kong, L. Li, and C. Shen, “SOLOv2: Dynamic and fast instance segmentation,” in *Adv. Neural Inf. Process. Syst.*, vol. 2020–December, 2020, Conference paper. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85107865426&partnerID=40&md5=e9e4f881a791ddbd81f564f74a345507>
- [123] Z. Wang, W. Liu, Y. Wang, and B. Liu, “Agcyclegan: Attention-guided cyclegan for single underwater image restoration,” in *ICASSP 2022–2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 2779–2783.
- [124] B. Xu, H. Liang, R. Liang, and P. Chen, “Locate globally, segment locally: A progressive architecture with knowledge review network for salient object detection,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 4, 2021, pp. 3004–3012.
- [125] S. Xu, M. Zhang, W. Song, H. Mei, Q. He, and A. Liotta, “A systematic review and analysis of deep learning-based underwater object detection,” *Neurocomputing*, vol. 527, pp. 204–232, 2023.
- [126] W. Yan, Y. Zhang, P. Abbeel, and A. Srinivas, “Videogpt: Video generation using vq-vae and transformers,” *arXiv preprint arXiv:2104.10157*, 2021.
- [127] C.-Y. Yang, H.-W. Huang, Z. Jiang, H. Wang, F. Wallace, and J.-N. Hwang, “A density-guided temporal attention transformer for indiscernible object counting in underwater videos,” in *ICASSP 2024–2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 5075–5079.
- [128] L. Yang, Z. Xu, H. Zeng, N. Sun, B. Wu, W. Cheng, J. Bo, L. Li, Y. Dong, and S. He, “Fishdb: an integrated functional genomics database for fishes,” *BMC Genomics*, vol. 21, 11 2020.
- [129] Z. Yang, L. Li, K. Lin, J. Wang, C.-C. Lin, Z. Liu, and L. Wang, “The dawn of lmms: Preliminary explorations with gpt-4v (ision),” *arXiv preprint arXiv:2309.17421*, vol. 9, no. 1, p. 1, 2023.
- [130] C.-H. Yeh, C.-H. Lin, L.-W. Kang, C.-H. Huang, M.-H. Lin, C.-Y. Chang, and C.-C. Wang, “Lightweight deep neural network for joint learning of underwater object detection and color conversion,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 11, pp. 6129–6143, 2021.
- [131] M. Zand, H. Damirchi, A. Farley, M. Molahasanji, M. Greenspan, and A. Etemad, “Multiscale crowd counting and localization by multitask point supervision,” in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 1820–1824.
- [132] C. Zhang, R. Cong, Q. Lin, L. Ma, F. Li, Y. Zhao, and S. Kwong, “Cross-modality discrepant interaction network for rgb-d salient object detection,” in *Proceedings of the 29th ACM international conference on multimedia*, 2021, pp. 2094–2102.
- [133] C. Zhang, L. Liu, G. Huang, H. Wen, X. Zhou, and Y. Wang, “Webuot-1m: Advancing deep underwater object tracking with a million-scale benchmark,” *arXiv preprint arXiv:2405.19818*, 2024.
- [134] J. Zhang, J. Huang, S. Jin, and S. Lu, “Vision-language models for vision tasks: A survey,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [135] P. Zhang, T. Yan, Y. Liu, and H. Lu, “Fantastic animals and where to find them: Segment any marine animal with dual sam,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 2578–2587.
- [136] W. Zhang, G. Chen, P. Zhuang, W. Zhao, and L. Zhou, “Catnet: Cascaded attention transformer network for marine species image classification,” *Expert Systems with Applications*, p. 124932, 2024.
- [137] Z. Zhao, C. Xia, C. Xie, and J. Li, “Complementary trilateral decoder for fast and accurate salient object detection,” in *Proceedings of the 29th acm international conference on multimedia*, 2021, pp. 4967–4975.
- [138] Z. Zheng, Y. Chen, J. Zhang, T.-A. Vu, H. Zeng, Y. H. W. Tim, and S.-K. Yeung, “Exploring boundary of gpt-4v on marine analysis: A preliminary case study,” *arXiv preprint arXiv:2401.02147*, 2024.
- [139] Z. Zheng, H. Liang, B.-S. Hua, Y. H. Wong, P. Ang, A. P. Y. Chui, and S.-K. Yeung, “Coralscop: Segment any coral image on this planet,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 28 170–28 180.
- [140] Z. Zheng, J. Zhang, T.-A. Vu, S. Diao, Y. H. W. Tim, and S.-K. Yeung, “Marinegpt: Unlocking secrets of ocean to the public,” *arXiv preprint arXiv:2310.13596*, 2023.
- [141] J. Zhong, M. Li, H. Zhang, and J. Qin, “Combining photogrammetric computer vision and semantic segmentation for fine-grained understanding of coral reef growth under climate change,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 186–195.
- [142] Y. Zhong, J. Yang, P. Zhang, C. Li, N. Codella, L. H. Li, L. Zhou, X. Dai, L. Yuan, Y. Li et al., “Regionclip: Region-based language-image pretraining,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 793–16 803.
- [143] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, “Learning to prompt for vision-language models,” *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.
- [144] Q. Zhou, S. Wang, J. Liu, X. Hu, Y. Liu, Y. He, X. He, and X. Wu, “Geological evolution of offshore pollution and its long-term potential impacts on marine ecosystems,” *Geoscience Frontiers*, vol. 13, no. 5, p. 101427, 2022.
- [145] C. Zhu, R. H. Byrd, P. Lu, and J. Nocedal, “Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale bound-constrained optimization,” *ACM Transactions on mathematical software (TOMS)*, vol. 23, no. 4, pp. 550–560, 1997.
- [146] P. Zhu, H. Wang, and V. Saligrama, “Don’t even look once: Synthesizing features for zero-shot detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 693–11 702.
- [147] P. Zhuang, Y. Wang, and Y. Qiao, “Wildfish++: A comprehensive fish benchmark for multimedia research,” *IEEE Transactions on Multimedia*, vol. 23, pp. 3603–3617, 2021.
- [148] Z. Ziqiang, X. Yaofeng, L. Haixin, Y. Zhibin, and S.-K. Yeung, “Coralvos: Dataset and benchmark for coral video segmentation,” *arXiv preprint arXiv:2310.01946*, 2023.
- [149] Z. Ziqiang, C. Yiwe, Z. Huimin, V. Tuan-Anh, H. Binh-Son, and Y. Sai-Kit, “Marineinst: A foundation model for marine image analysis with instance visual description,” *European Conference on Computer Vision (ECCV)*, 2024.