MARINEMAID: DATASET AND BENCHMARK ON DE-TECTING AND UNDERSTANDING MARINE CREATURES

Anonymous authors

000

001

003 004

006

007 008 009

010

018

019

021

023

025

026

027

028 029

031

032

034

039

040

041

042

043

044

045

046

048

Paper under double-blind review



Figure 1: We present MarineMaid dataset, the first dataset and benchmark specifically designed for marine visual understanding. *Top:* MarineMaid consists of two main components: high-quality bounding boxes for object detection and fine-grained instance captions for marine vision-language understanding. *Bottom:* MarineMaid enables a wide range of marine visual understanding tasks, including *object detection, open-vocabulary object detection, instance captioning* and grounding.

ABSTRACT

Oceans, covering more than 70% surfaces of our blue planets are less explored by the whole computer vision community. The scarcity of the labeled data is attributed to the most hindering issue. In this work, we propose a novel and comprehensive dataset called MarineMaid specifically designed for marine monitoring and understanding, including a wide spectrum of marine creatures. Based on the essential requirements of the marine research community, we adopt *object detection* and vision-language understanding as our two fundamental tasks. The former object detection could yield precise localization and category predictions for species identification and monitoring. Besides the sole category and BBOX predictions, the latter vision-language understanding generates redundant and comprehensive captions about biological traits required for domain experts. MarineMaid contains 12,873 fine-grained instance-captioning pairs and 42,217 bounding boxes annotated by domain experts. We have benchmarked 14 state-of-the-art algorithms on our MarineMaid dataset to reveal the strengths and limitations of existing generalpurpose and domain-specific algorithms. The hierarchical and comprehensive experimental results provide valuable insights on how to develop practical and efficient marine visual perception algorithms to satisfy the domain requirements. To foster the further development of this direction, we will release our MarineMaid dataset with the acceptance of this paper.

050 1 INTRODUCTION

The unbounded depths of the ocean (Epstein et al., 1993; Ormond et al., 1997), rich with mysteries,
 have driven researchers to explore relentlessly, aiming to uncover its hidden secrets and valuable
 treasures (Thorne-Miller, 1999). The marine ecosystem (Epstein et al., 1993; Halpern et al., 2008)

054 is the most productive of all underwater ecosystems and shares immense ecological, social, and 055 economic value. Performing marine study plays a significant role in protecting the marine environment 056 and understanding marine science. However, marine research is limited compared with its volume. 057 The ocean is continuously being polluted, leading to the migration of marine organisms (Perry et al., 058 2005) and species changes (Hiddink & Ter Hofstede, 2008; Poloczanska et al., 2013). Automatic marine life detection algorithms based on computer vision techniques are keenly required. Existing methods for monitoring and assessing marine ecosystem changes suffer from inefficiency and high 060 labor costs. Nevertheless, utilizing computer vision techniques and deep learning algorithms can 061 rapidly analyze marine images and videos, identifying species (Khan et al., 2023b) and tracking their 062 migrations (Danovaro et al., 2010). 063

064 Object detection (Redmon et al., 2016; Ren et al., 2015; 2016; Liu et al., 2016), as a fundamental task, is to localize the interests of objects while discriminating the category information. The detected 065 objects with bounding box annotations are important for species identification (Khan et al., 2023b), 066 object tracking (Alawode et al., 2022), and object counting (Sun et al., 2023). To boost efficient 067 marine object detection, there are several efforts proposed to build datasets (Zhuang et al., 2020; 068 Liu et al., 2021b) and benchmarks (urp, 2020) to optimize powerful object detection algorithms. 069 However, the categories (e.g., sea urchin, shark and etc.) are very limited, which are far away from satisfying to monitor a large range of marine creatures. Furthermore, the category information is 071 not sufficient to satisfy the monitoring and surveying requirement, where biological traits (Miatta 072 et al., 2021; Costello et al., 2015) are usually required. Furthermore, due to the essential monitoring 073 purpose, the algorithms should also be able to detect a wide spectrum of objects (Zheng et al., 2023) 074 and demonstrate strong generalization ability to unseen marine objects.

075 Vision-language models (VLMs) (Liu et al., 2024; 202, 2023; Team et al., 2023) achieve remarkable 076 success thanks to large-scale datasets (Krishna et al., 2017; Gurari et al., 2018; Kazemzadeh et al., 077 2014; Shao et al., 2019). The VLMs could yield redundant visual descriptions based on the visual 078 inputs, describing the objects with detailed attributes (e.g., color, pose, activity, and etc). Despite the 079 remarkable success of VLMs in a large number of visual understanding tasks, they are still poorly known to generate reasonable and domain-specific visual understanding for marine creatures. There 081 are two main limitations when directly utilizing existing VLMs for marine visual understanding: data distribution shift and the lack of ability to localize and then describe the biological traits of the marine instances. The existing VLMs were mainly driven by datasets with dominant in-air objects 083 and very limited marine objects, leading to unsatisfactory marine object understanding ability. There 084 is a gap in the development and evaluation of VLMs for marine visual understanding for scientific 085 research purposes. Furthermore, VLMs are optimized by redundant image-text pairs that succeed in holistic view understanding but struggle with detecting and understanding specific marine creatures 087 with irregular boundaries/poses and also the ability to camouflage themselves into the background. 088 Besides, generating the biological traits for detected instances with detailed descriptions of the spatial 089 information/relationship between objects is also important to yield a complete analysis report. There 090 is still a gap in utilizing existing algorithms and datasets for domain-specific marine research.

091 To fill this gap, we propose the first marine dataset and benchmark called MarineMaid to achieve 092 robust and accurate marine visual understanding with detailed descriptions of biological traits from 093 various aspects. Our dataset with rich biodiversity comprises 14,645 marine images with more 094 than 42k human-labeled bounding boxes and instance captions, enhancing the understanding of the complex marine ecosystems. MarineMaid dataset enables various tasks, including open-vocabulary 096 object detection, region-specific image/instance captioning and visual grounding specifically designed for marine creatures. Unlike existing image-text datasets with only short descriptions, MarineMaid 098 provides a comprehensive and detailed description (average word length is 42) of the biological traits of the marine creatures from 4 aspects. To the best of our knowledge, our MarineMaid dataset is the 099 first marine dataset to support marine monitoring and further analysis. 100

Based on our MarineMaid dataset, we have benchmarked the existing object detection, VLMs, and grounding algorithms to explore the boundary of these advanced algorithms to perform detailed marine visual understanding. Our MarineMaid stands as a novel and challenging testbed for both computer vision and marine research communities. Our main contributions are as follows:

- 105
- 106 107
- We propose the first region-level instance-caption pair dataset specifically designed for marine creatures, containing 12,873 fine-grained instance-captioning pairs and 42,217 BBOXs annotated by domain experts.

• We benchmark various marine visual understanding tasks including close-set object detection, open-vocabulary object detection, visual grounding, and instance captioning based on 14 state-of-the-art models.

112 2 RELATED WORK

108

110

111

113 Existing Marine Research. Unlike our everyday stuff, marine creatures usually possess significant 114 diversity (a wide spectrum of poses, appearance, and patterns). Performing efficient marine visual 115 understanding could harness the advanced algorithms (Li et al., 2021; Hong et al., 2020) to elevate 116 marine research, conservation, and industrial endeavors. Existing marine datasets (e.g., MAS3K (Li 117 et al., 2020; 2021), WildFish (Zhuang et al., 2018), WildFish++ (Zhuang et al., 2020), SUIM (Islam 118 et al., 2020)) have been proposed for promoting the recognition performance of marine organisms. 119 However, most of these datasets only contain a few pre-defined categories without detailed captions, 120 which limits the ability to accelerate the accumulation of detailed marine visual analysis through 121 the creation of ocean databases and scientific data. Meanwhile, domain knowledge and expertise 122 are required to do high-quality annotations (for both BBOX and caption), which is costly and timeconsuming. In this work, we aim to propose a large-scale marine dataset with a wide spectrum of 123 marine creatures. 124

125 **Object Detection**. Object detection is a fundamental computer vision problem (Lin et al., 2014; 126 Ren et al., 2015; 2016), localizing the interests of objects and discriminating object categories 127 simultaneously. The detection algorithms mainly fall into two categories: 1) one-stage algorithms (Liu et al., 2016; Ge et al., 2021; Redmon et al., 2016) perform localization and classification in parallel; 2) 128 two-stage detection algorithms (Ren et al., 2015; 2016; He et al., 2017) generate the object proposals 129 and then perform localization regression. However, these algorithms mainly perform close-set object 130 detection. To address this limitation, open-vocabulary object detection (OVOD) (Zareian et al., 2021; 131 Yao et al., 2023; Kim et al., 2023; Wang et al., 2023) aims to generalize beyond the limited number 132 of pre-fixed classes during the training phase. The goal is to detect novel classes at the inference 133 stage. The dominant way of performing OVD is to adopt a pre-trained visual encoder from a trained 134 cross-modality alignment model, which is optimized by millions of image-text pairs from public 135 websites. RegionCLIP (Kim et al., 2023) proposed to perform the regional visual feature and the 136 textual conception alignment to promote the generalization ability to *unseen* categories.

137 Vision-Language Understanding. Vision-language models (VLMs) (Liu et al., 2024; 202, 2023; 138 Team et al., 2023; Zhu et al., 2023; Liu et al., 2023a; Zheng et al., 2023; Li et al., 2022; 2023a) achieve 139 remarkable success thanks to large-scale datasets such as Visual Genome (Krishna et al., 2017), 140 VizWiz (Gurari et al., 2018), RefCOCO (Kazemzadeh et al., 2014), and Objects365 (Shao et al., 2019). 141 VLMs bridge vision modality and text modality together to harness the power of large language 142 models (LLMs) (OpenAI, 2022; 2023) and vision encoders (Dosovitskiy et al., 2020). Optimized by 143 millions of image-text pairs, CLIP (Radford et al., 2021) demonstrated a strong zero-shot recognition ability for diverse images. BLIP (Li et al., 2022; 2023a) bootstraps vision-language pre-training 144 from frozen pre-trained image encoders and frozen language decoders. However, these datasets 145 only contain in-air objects or very limited marine objects, which is due to the poor ability of marine 146 domain tasks. Furthermore, VLMs also struggle with the limited ability to perform region-level 147 instance understanding following the user instructions. 148

¹⁴⁹ 3 DATASET AND APPROACH

Overview. We start by elaborating on the detailed dataset construction procedure of our MarineMaid
 and outlining the characteristics of our dataset, along with relevant statistics and explanations. We
 then provide the hierarchical and extensive experiments to benchmark marine object detection (in cluding both close-set and open-vocabulary formulations), visual grounding, and instance captioning,
 revealing the strengths and limitations of existing algorithms.

156 3.1 DATASET CONSTRUCTION

157 Data collection. We collect images from the Internet. To maintain data quality and diversity, we
 158 manually reviewed all the images and removed duplicates or instances that did not align with the pre 159 defined categories. Existing datasets (Schuhmann et al., 2021) mainly utilized alt-texts to formulate
 160 the image-text pairs (*image-level*). However, the texts suffer from limited information (short captions),
 161 misalignment with the visual contents, and deviation from domain-specific requirements. In contrast,
 we generate comprehensive and contextually relevant *instance* captions based on the domain experts.



Figure 2: Overview of the dataset construction and data labeling pipeline, which can be summarized 177 into five stages: 1) marine object categories are generated based on ChatGPT-3.5/GPT-4; 2) crawling 178 corresponding marine images from the Internet (mainly from Google image engine and Flickr); 3) 179 we employ SAM (Kirillov et al., 2023) model to label all marine objects present in each image by 180 receiving the human prompts to iteratively obtain high-quality BBOX annotations; 4) the cropped 181 image region based on the annotated BBOX from the whole image serves as the input to domain-182 specific VLM (MarineGPT (Zheng et al., 2023) used in this work) to generate object instance caption candidates for further human refinement; 5) domain experts refine the generated captions 183 from some pre-defined aspects as the positive instance captions. We also provide the additional 184 binary annotations from 11 diverse properties to identify some common prediction errors produced 185 by VLMs and we regard these captions with binary annotations as negative captions (discussed in Supplementary material). 187

To promote labeling efficiency, we first utilize the marine-specific VLM MarineGPT to generate the
 candidate captions and the domain experts perform the refinement and revision from pre-defined
 aspects.

191 Specific features: 1) Wide spectrum of marine object categories (670 categories), varying from 192 Cephalopods, Crustaceans, Sharks, Rays, Reptiles, Mammals, Aves, Corals, to Invertebrates. 2) 193 **Hierarchical taxonomy**: including 6 coarse-to-fine granularities (Kingdom, Phylum, Class, Order, 194 Family, Genus) by automatically querying the official Worms (Ahyong et al., 2024) database. 3) 195 **Image diversity**: images were captured in various environmental conditions (*e.g.*, deep-sea, blurs, 196 clutters, aquariums, markets, *etc*). Meanwhile, the images describe object instances from different aspects: activity events (e.g., hunting, reproductive, interactive, etc), life stages (e.g., juvenile and 197 adult), and image styles. 198

Positive vs. Negative.We also provide the *additional information* for the negatives to reveal common mistakes made by the models. We define 11 properties: classification, background, unexisting, spatial, action, size, color, shape, texture, material, and counting to summarize wrong captions. These negatives from VLMs and human post-processing offer more valuable insights compared to (Zhao et al., 2022; Yuksekgonul et al., 2022) that replaced correct objects with random noun phrases. These negative samples serve as the hard negatives to force the model to learn and recognize subtle feature differences.

Data statistics. Our dataset consists of 14,645 images, from a total of 670 categories. We manually
labeled all identifiable marine life object, resulting in a total of 42,217 labeled bounding boxes.
Among these, there are 24,197 large, 10,555 medium, and 7,465 small bounding boxes. There are
12,873 captions that have been refined by domain experts specializing in marine specialties, resulting
in a superior level of quality. The average length of these refined positive captions is 42. Totally we
have 22,321 refined and generated positive captions, and 12,431 generated negative captions.

2122133.2LABELING PIPELINE

Our labeling pipeline encompasses three main stages: 1) BBOX labeling and refinement; 2) caption
 generation from VLM and refinement based on domain experts; and 3) cross-checking verification.
 BBOX generation. We first manually label bounding boxes for all the recognizable marine organisms



Figure 3: We provide the data statistics and the class distribution of MarineMaid at the Class-level granularity. The *seen* and *unseen* classes are split to perform open-vocabulary object detection.

within the image to perform dense labeling. To ensure the quality of the labeled BBOX annotations, 232 we perform further refinement to ensure the whole instance (e.g., the transparent tail of the fish, and 233 the slender legs of the shrimp) is accurately labeled. Caption generation and refinement. We first 234 crop the marine object instance based on the BBOX annotations and feed the cropped image region 235 to MarineGPT to generate the caption candidate based on the prompt "describe this image in detail". 236 Please note that we only generate descriptive and informative captions based on image regions larger 237 than 1024 pixels. Then based on the caption candidate, the domain experts do the refinement from 238 four aspects: features (e.g., unique characteristics, injuries, color, shape, size, etc), spatial information 239 (e.g., absolute and relative position), background and activity events (e.g., individual or mutual). For 240 each image, we only select one to perform caption refinement and we provide additional tags from the 241 pre-defined 11 properties for the incorrect captions to formulate the hard negatives. Cross-checking 242 validation. Finally, we perform the cross-validation based on two annotators to revise potential errors. 243 Filtering: filters out some model-generated prompts or "unrecognizable" content. The annotators will cross-check and correct evident errors in captions, tags, and bounding boxes. Experts finally conduct 244 inspections and verifications on uncertain objects to ensure accuracy and reliability. The construction 245 of our MarineMaid dataset involves 16 domain experts with 624 human hours in total. 246

247 3.3 COMPARISON WITH EXISTING DATASETS AND BENCHMARKS

248 We provide a direct comparison with existing general-purpose and domain-specific datasets in Table 1 249 from various aspects: the data/annotation volume; annotation type; whether the image/instance 250 captions provided and the average word length of these corresponding captions; and Taxonomy for 251 hierarchical classification and understanding. MarineMaid dataset possesses three main advantages over existing datasets: 1) compared with existing marine datasets, which mainly provide the BBOX 253 and mask annotations, the MarineMaid dataset provides detailed instance captions for the object instance besides the BBOX annotations. 2) Compared with general-purpose datasets that contain a 254 large scale of image/instance captions, our provided instance captions are significantly longer (42 vs. 255 12), describing diverse biological traits of marine creatures. 3) Compared with the existing Wildfish++ 256 dataset with both taxonomy and visual descriptions from the domain experts, MarineMaid is 10 times 257 larger and contains a wide range of marine creatures while wildfish++ only focuses on fish. 258

4 EXPERIMENTS

259

229

230 231

In this section, to comprehensively evaluate the effectiveness of marine visual understanding, we
 choose three representative visual understanding tasks, including object detection (both close-set and
 open-vocabulary settings), region-level instance captioning, and visual grounding. We benchmark the
 existing state-of-the-art algorithms for corresponding tasks on our MarineMaid dataset.

265 4.1 Object Detection

Experimental settings. Dataset split. We construct *seen/unseen* split following three settings:
1) Class-level: We consolidated the 670 categories into 33 categories based on their taxonomic Class (Ahyong et al., 2024). Certain object categories (*e.g.*,, bryozoa), are classified at a higher level of granularity (Phylum) and are therefore excluded from the Class-level categories. As illustrated in Fig. 3, we adopt 24 Classes as *seen* and the other 9 Classes as *unseen*. 2) Intra-Class: Intra-Class

Table 1: We provide a direct comparison between our MarineMaid dataset with both general-purpose 270 datasets and marine-specific datasets. - indicates that the numbers were either not reported in their 271 publications or we are unable to conduct statistical analysis. 272

273	Datasets	Categories	Images	Annotations	Image/Instance Captions	Avg. Length	Taxonomy
274	DUO (Liu et al., 2021a)	4	7,782	74,515 _{bbox}	None	None	×
275	SUIM (Islam et al., 2020)	8	1,525	$1,525_{mask}$	None	None	×
	MAS3K (Li et al., 2020)	37	3,103	$3,103_{mask}$	None	None	×
276	UIIS (Lian et al., 2023)	7	4,628	$4,628_{mask}$	None	None	×
277	SEAMPD21 (Boulais et al., 2021)	130	28,328	$90,000_{bbox}$	None	None	×
211	Wildfish (Zhuang et al., 2018)	1,000	54,459	54,459 _{cls}	None	None	×
278	FishNet (Khan et al., 2023a)	17,357	94,532	$114,375_{bbox}$	None	None	\checkmark
279	Wildfish++ (Zhuang et al., 2020)	2,348	103,034	$103,034_{cls}$	3,187	56	√
000	nocaps (Agrawal et al., 2019)	-	15,100	Caption	166,100	-	×
200	Redcaps (Desai et al., 2021)	_	12,011,121	Caption	12,011,121	9	×
281	Pascal Sentences (Rashtchian et al., 2010)	20	1,000	Caption	4,998	10	×
202	SBU Captions (Ordonez et al., 2011)	81	1,000,000	Caption	1,000,000	12	×
202	MarinaMaid	670	14 645	42 217	12,873 (Refined)	42	1
283		070	14,045	42,217bbox	34,752 (All)	33	✓

284

287

289

290

categorization is obtained by retrieving object categories at the Class-level. Under this setting, we 286 have 555 seen categories and 109 unseen categories. 3) Inter-Class: we choose 1 object category from every 4 object categories in each Class as the *unseen* and the other 3 object categories as *seen*. 288 We omit the Class that contains less than 4 object categories. With this setup, there are 482 seen categories and 161 unseen categories.

Implementation details. Close-set object detection setting. We mainly include 3 close-set object 291 detection algorithms (Faster-RCNN (Ren et al., 2015), GridRCNN (Lu et al., 2019) and YOLOX (Ge 292 et al., 2021)) and report the mAP₅₀ of 24 seen categories under three settings(Class-level, Intra-293 Class and Inter-Class). Our implementation of these models is based on MMDetection (Chen et al., 294 2019) using the official experimental setting. Please note that we do not evaluate these close-set 295 object detection algorithms on the *unseen* categories. **Open-Vocabulary Object Detection** We 296 evaluate the performance of 3 open-vocabulary object detection algorithms (RegionCLIP (Zhong 297 et al., 2022), UniDetector (Wang et al., 2023) and DECOLA (Cho & Krähenbühl, 2023)) on our 298 MarineMaid dataset. For RegionCLIP (Zhong et al., 2022), we follow the official experimental 299 setting and fine-tune the model on our MarineMaid dataset. We adopt the single-dataset training strategy for UniDetector (Wang et al., 2023) to continuously optimize it in an end-to-end fashion. For 300 DECOLA (Cho & Krähenbühl, 2023), we utilize their best-performing model with Swin-B backbone 301 (phase 1) as the pre-trained model. We inherit the language-conditioned detection training procedure 302 of DECOLA while keeping other configurations the same. At the evaluation stage, we report the 303 quantitative results for both seen and unseen categories. The mAP₅₀ is computed to comprehensively 304 evaluate the ability of models to detect overall marine object instances. 305

306 **Comparison and analysis.** We report the quantitative result in Table 2 and all the experiments are conducted following the same train/val data split. We have observed that the existing generalist 307 object detection algorithms still face challenges when optimized by underwater images in providing 308 accurate object localization. This can be attributed to two potential reasons: 1) the huge conception 309 distance between the in-air object categories and marine object categories; and 2) the diversity of 310 underwater data and the inherent challenges of underwater scenes make it difficult to extract features. 311 Furthermore, as demonstrated, open-vocabulary detection algorithms, continuously fine-tuned on 312 the MarineMaid dataset, typically exhibit improved detection performance even on *seen* categories 313 compared to close-set counterparts. We attribute such promoted performance to the optimization 314 through large-scale datasets with redundant supervised training data during the pre-training procedure. 315 We present a qualitative comparison of the results in Fig. 4 under Class-level setting. DECOLA 316 exhibits superior performance in semantic and object localization when detecting *seen* objects. However, when it comes to *unseen* objects, the models struggle to accurately classify the object 317 category. In both Intra-Class and Inter-Class settings, DECOLA is the sole model to gain an advantage 318 over the fine-grained marine species. We attribute such powerful fine-grained recognition ability to 319 its language-conditioned query selection strategy. 320

- 321 4.2 INSTANCE CAPTIONING 322
- Experimental settings. We benchmark off-the-shelf VLMs from two aspects: *image-level* and 323 region-level. The former image-level VLMs (LLAVA (Liu et al., 2024), MiniGPT-4 (Zhu et al.,



Figure 4: The qualitative comparison between different algorithms under the Class-level setting. *The left part of the dashed line*: the results of *seen* category. *The right part*: the results of *unseen* category.

Table 2: Quantitative object detection (close-set and open-vocabulary) results on our MarineMaid dataset. – indicates the results cannot be computed under the settings.

Method	Class-level	Seen Intra-Class	Inter-Class	Class-level	Unseen Intra-Class	Inter-Class
FasterRCNN (Ren et al., 2015) YOLOX (Ge et al., 2021) GridRCNN (Lu et al., 2019)	28.7 27.5 32.7	17.6 21.7 28.1	16.7 21.0 28.6	- - -	- -	- - -
UniDetector (Wang et al., 2023) RegionCLIP (Zhong et al., 2022) DECOLA (Cho & Krähenbühl, 2023)	31.5 <u>39.8</u> 66.7	23.3 <u>34.1</u> 88.8	24.1 <u>29.8</u> 86.9	8.2 <u>12.2</u> 37.7	0.4 <u>6.2</u> 51.6	$ \frac{0.7}{0.4} $ 52.3

2023), BLIP2 (Li et al., 2023b) and InstructBLIP (Dai et al., 2024)) were optimized by image-level captions and lacked the ability to understand specific object instances. We evaluate these image-level VLMs based on the following user instruction: "describe the object in this figure". The latter region-level VLMs (GroundingLMM (Rasheed et al., 2023), GPT4RoI (Zhang et al., 2023)) were optimized by paired image region prompts and the corresponding instance captions. We provide the BBOX annotation in the given text prompt following the experimental setting of (Rasheed et al., 2023; Zhang et al., 2023). We perform the evaluations based on the positive instance captions to analyze their capability in describing marine instance objects. To quantitatively measure the performance of various algorithms, we adopt the widely used captioning metrics (Hessel et al., 2021; Vedantam et al., 2015; Banerjee & Lavie, 2005; Lin, 2004; Papineni et al., 2002) (including CLIPScore, RefCLIPScore (Hessel et al., 2021), CIDEr (Vedantam et al., 2015), BLUE-4 (Papineni et al., 2002), METEOR (Banerjee & Lavie, 2005) and Rouge (Lin, 2004)) to compute quantitative results in Table 3. Besides the human-constructed instance captions proposed in our MarineMaid dataset, we also construct a starting sentence to include the category information for the selected object instance: "This is a <*Category Name>*.", where the <*Category Name>* is the placeholder to compensate the scientific category-level information of each instance. In this way, by penalizing generated plausible but not domain-specific responses (e.g., "fish", "animal" and "mammal"), we encourage the model to generate the scientific captions to satisfy the domain requirements.

Implementation details. We perform the evaluation only based on the released official models
 provided by various algorithms on our MarineMaid dataset and our experiments were conducted using
 an NVIDIA L20 GPU. For LLAVA (Liu et al., 2024), we choose its V1.5-7b version for evaluation.

Table 3: Results of various algorithms (image-level and region-level) on instance captioning.

373	Method	CLIPScore↑	RefCLIPScore↑	CIDEr↑	BLUE-4↑	METEOR↑	Rouge↑
374 375 376	LLAVA (Liu et al., 2024) MiniGPT-4 (Zhu et al., 2023) BLIP2 (Li et al., 2023b) InstructBLIP (Dai et al., 2024)	73.78 74.48 <u>76.22</u> 76.60	72.27 73.43 <u>73.73</u> 75.25	4.93 5.72 <u>9.96</u> 12.09	8.77 7.18 8.16 13.94	<u>7.70</u> 16.90 5.95 7.40	20.76 28.03 18.96 <u>21.31</u>
377	GroundingLMM (RefCOCOg) (Rasheed et al., 2023) GroundingLMM (VG) (Rasheed et al., 2023) GPT4RoI (Zhang et al., 2023)	73.04 71.15 71.28	70.97 69.04 68.71	4.37 4.06 3.53	4.39 2.47 2.81	4.60 4.11 4.07	16.37 15.22 15.08

3				
)				
		The state of the		
LLAVA	A black fish with a white stripe on its tail swims through the water.	A person is holding a fish in their hand, showing off its shiny scales.	A small crab is sitting inside a pink <mark>shell,</mark> possibly a sea anemone.	
	this is a black fish with white stripes swimming in a tank filled with water, the fish has a long	this image shows a group of people holding a large fish on a boat, the fish is a blue marlin	this image shows a large, pink sea anemone with its tentacles spread out, surrounded by	
Minicott	slender body and a small mouth with sharp teeth. it is swimming towards the camera, its fins flapping rapidly as it moves through the water, the background of the image is a dark, racky aging with a small amount of alcan	which is a large and powerful fish that is prized for its meat and sport fishing. the people in the	other sea creatures such as a crab and a fish the anemone appears to be in a coral reef	
MinigP1-4		image are wearing fishing gear and appear to be preparing to release the fish back into the	environment, with other coral and seaweed visible in the background. the image is taken	
	rocky cave with a small amount of algae growing on the walls.	water.	the water to create a beautiful, colorful scene.	
BLIP2	a black slender <mark>ray</mark> swimming in an aquarium	a tuna caught on a fly fishing trip in the bahamas	a close up of the inside of a sea urchin	
Instruct BLIP	a black and white fish swimming in an aquarium	a man holding a large fish on a boat	a close up of a large pink <mark>shell</mark> on the ocean floor	
Groundinរួ LMM	a black fish with a white stripe on its head	a large fish being held by two men	a purple and pink <mark>sea sta</mark> r	
GPT4Rol	a black fish swimming in the water	a man holding a fish	a purple and pink organism	
	This is a black ghost knifefish swimming in the	The object in the middle is a tuna being held in	This is a close-up of a pink barrel sponge seen	
GT	water. The water is clear and there is a piece of rock, aquatic plants and gravel in the	the left hand of a fisherman in a boat. The fish is silver with purple and green iridescent markings	from above the animal. The inside surface is smooth and layered and the outside surface	
5	background. The body of the fish is black and it has a long tail with black and white stripes.	on its body. Its dorsal fins are separate and both steep and pointed. There is another person standing to the fisherman's left.	appears spikey and rough. The appears to be a smaller blue sponge behind the object.	

Figure 5: The qualitative results of different algorithms on marine object instance understanding. Best viewed in color.

400 401

398

399

The language model of MiniGPT-4 (Zhu et al., 2023) is set to LLaMA-2 (Touvron et al., 2023).
As for the GroundingLMM (Rasheed et al., 2023), we report the results of the models fine-tuned on RefCOCOg dataset (Kazemzadeh et al., 2014) and Visual Genome (VG) dataset (Krishna et al., 2017), respectively.

406 **Comparison and analysis.** All the quantitative results are reported in Table 3. Please note that 407 CLIPScore and RefCLIPScore (Hessel et al., 2021) are computed based on the whole image. We observe that image-level VLMs achieve various scores when there are human-constructed reference 408 captions. LLAVA (Liu et al., 2024) and BLIP2 (Li et al., 2023b) achieve very poor outputs on 409 CIDEr (Vedantam et al., 2015) and BLUE-4 (Papineni et al., 2002) since these two tend to generate 410 very short answers and they also make some wrong recognitions. InstructBLIP (Dai et al., 2024) 411 performs best on CLIPScore, RefCLIPScore (Hessel et al., 2021), CIDEr (Vedantam et al., 2015) 412 and BLUE-4 (Papineni et al., 2002). This indicates that the instruction-following tuning could 413 heavily promote the ability of the models to understand the instances following the user instructions. 414 However, the generated instance captions are still too short to satisfy the domain requirement. Region-415 level VLMs also achieve very poor results since they cannot accurately localize the specific marine 416 instances by the user-provided BBOX prompts. Thus, the region-level VLMs still describe the whole 417 image and yield wrong captions as demonstrated in Fig. 5. There is still a gap when utilizing existing 418 VLMs for marine instance understanding.

419 420

421

4.3 GROUNDING

422 **Experimental settings.** We finally demonstrate the performance of existing general-purpose ground-423 ing algorithms in handling marine visual localization. Using our captions as prompts, we apply 424 these algorithms to generate target bounding boxes, which are then compared to the ground truth 425 bounding boxes in our dataset. Specifically, we examine GroundingDINO (Liu et al., 2023b) and 426 GroundVIP (Shen et al., 2023). These models are not fine-tuned on our training set and are directly 427 evaluated on our validation set. Additionally, captions that are negatives, empty, and with no noun 428 phrases detected by nltk package are excluded to ensure a smooth evaluation process. The results are 429 reported in Table 4, following the default evaluation metrics (Recall used in GroundingDINO (Liu et al., 2023b) and accuracy for GroundVLP (Shen et al., 2023)). To guarantee a proper configuration 430 of the evaluation environment settings, we meticulously adhered to the instructions and evaluation 431 program provided by the authors (discussed in Supplementary).

(d) Query: This is a sea urch (a)Query: This is a bubble-tip (b) Query: This is an angel shark (c) Query: This is a shark in a tank n with long under aquarium lighting. It has swimming on the ocean floor. It is large in size with a flat body. Its tail It appears to be swimming in the sharp black spines on its body. It is located in a body of water with bulbous tops on its tentacles that water are closely packed together in is long with two raised dorsal fin other sea creatures nearby. There is purple color, with a large protruding The color of its body is similar to the a purple and yellow fish on top of it spot at the ocean floor.

Figure 6: Results of GroundingDINO (Liu et al., 2023b) (a,b) and GroundVLP (Shen et al., 2023) (c,d) on our MarineMaid dataset. Green texts and BBOXs indicate the query and GT respectively. Red texts and BBOXs indicate model-generated predictions and corresponding BBOX outputs.

Table 4: Quantitative visual grounding results on MarineMaid dataset.

Method	Evaluation Metric	Class-level	Seen Intra-Class	Inter-Class	Class-level	Unseen Intra-Class	Inter-Class
GroundingDINO (Liu et al., 2023b)	R@1	38.8	37.5	37.2	46.8	46.8	45.6
GroundingDINO (Liu et al., 2023b)	R@5	67.3	65.5	66.2	76.9	76.8	74.5
GroundingDINO (Liu et al., 2023b)	R@10	78.0	76.0	76.8	84.9	84.9	83.3
GroundingVLP (Shen et al., 2023)	Accuracy	19.8	45.8	42.8	34.3	35.6	52.4

Implementation details. For GroundingDINO (Liu et al., 2023b), we employ the best Swin-B pretrained model (MM-GDINO-B*) as the backbone. The evaluation is conducted on a single GeForce RTX 2080 Ti. Other configurations are consistent with the original paper. For GroundVLP (Shen et al., 2023), we use the ALBEF and Swin-B Detic (Zhou et al., 2022) models provided by the authors, evaluated on our validation dataset.

Comparison and analysis. As reported in Table 4, there is an obvious performance drop (still 458 unsatisfactory performance) when utilizing the two grounding algorithms on marine creature local-459 ization. We attribute this to the gap in knowledge between everyday objects and marine creatures. 460 Furthermore, as depicted in Fig. 6, these algorithms struggle to accurately recognize and locate target 461 marine creatures, being hindered by the knowledge acquired from in-air objects. For instance, in 462 Fig. 6 (c), GroundVLP (Shen et al., 2023) mistakenly identifies a shark as a cow and fails to locate 463 the correct target based on the described action in the caption. Conversely, GroundDINO (Liu et al., 464 2023b) (b) correctly identifies the angle shark, but mistakenly recognizes its tail fin as dorsal fins, 465 further revealing the lack of ability to perform accurate marine visual grounding.

466 467 468

469 470

471

472

473

432

433

434

443

444

445

5 DISCUSSIONS AND CONCLUSION

New benchmark. The proposed MarineMaid serves as a novel comprehensive and diverse benchmark meticulously curated for marine research. Our dataset is introduced to enhance the assessment of existing algorithms for marine visual understanding. It includes a wide range of marine creatures across various environments, providing a valuable benchmark for testing and developing new models.

Broader impact. The study of marine creatures has several important applications, such as identifying and safeguarding rare animal species, preventing wildlife trafficking, and aiding in search-and-rescue operations. Our dataset deliberately excludes any military or sensitive scenes, ensuring its focus remains on benign and beneficial applications.

Limitation. Even though we tried our best to cover the most common marine creatures, we have to admit that the amounts of existing marine creatures are much larger than the included marine object categories. Our dataset will be continuously growing to include more marine object categories.

Conclusion. In this work, we propose the first large-scale marine datasets to enable both object
 detection and vision-language understanding. Our dataset supports various tasks, including *close- set object detection, open-vocabulary object detection, instance captioning*, and *grounding*. The
 comprehensive evaluation sheds light on the strengths and limitations of both general-purpose and
 domain-specific algorithms.

486 REFERENCES

488

489 490

491

492

496

Urpc dataset. https://openi.pcl.ac.cn/OpenOrcinus_orca/URPC2020_ dataset/datasets, 2020.

- Gpt-4v(ision) system card. 2023. URL https://api.semanticscholar.org/CorpusID: 263218031.
- Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra,
 Devi Parikh, Stefan Lee, and Peter Anderson. nocaps: novel object captioning at scale. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 8948–8957, 2019.
- S. Ahyong, C.B. Boyko, N. Bailly, J. Bernot, R. Bieler, S.N. Brandão, M. Daly, S. De Grave, S. Gofas, 497 F. Hernandez, L. Hughes, T.A. Neubauer, G. Paulay, B. Boydens, W. Decock, S. Dekeyzer, L. Van-498 depitte, B. Vanhoorne, R. Adlard, S. Agatha, K.J. Ahn, B. Alvarez, M.R.W. Amler, V. Amorim, 499 A. Anderberg, S. Andrés-Sánchez, Y. Ang, D. Antić, L.S.. Antonietto, C. Arango, T. Artois, 500 S. Atkinson, K. Auffenberg, B.G. Baldwin, R. Bank, A. Barber, I. Bartsch, D. Bellan-Santini, 501 N. Bergh, A. Berta, T.N. Bezerra, S. Blanco, I. Blasco-Costa, M. Blazewicz, L.A. Błedzki, P. Bock, M. Bonifacino, R. Böttger-Schnack, P. Bouchet, N. Boury-Esnault, R. Bouzan, G. Boxshall, R. Bray, A.L. Brito Seixas, N.L. Bruce, A. Bruneau, N. Budaeva, J. Bueno-Villegas, S. Cairns, J. Calvo Casas, P. Cárdenas, E. Carstens, P. Cartwright, T. Cedhagen, B.K. Chan, T.Y. Chan, 504 H. Choong, M. Christenhusz, M. Churchill, A.G. Collins, G.E. Collins, K. Collins, L. Consorti, 505 D. Copilas-Ciocianu, L. Corbari, R. Cordeiro, V.M.d.M. Costa, P.H. Costa Corgosinho, M. Coste, M.J. Costello, K.A. Crandall, F. Cremonte, T. Cribb, S. Cutmore, F. Dahdouh-Guebas, M. Daneliya, 507 J.C. Dauvin, P. Davie, C. De Broyer, P. de Lima Ferreira, V. de Mazancourt, L. de Moura Oliveira, H.A.B., de Sá, N.J. de Voogd, P. Decker, D. Defaye, H. Dekker, J.L. d'Hondt, I. Di Capua, S. Dippenaar, M. Dohrmann, J. Dolan, D. Domning, R. Downey, N. Dreyer, U. Eisendle, M. Eitel, 510 M. Eleaume, H. Enghoff, J. Epler, P. Esquete Garrote, N.L. Evenhuis, C. Ewers-Saucedo, M. Faber, 511 D. Figueroa, C. Fišer, E. Fordyce, W. Foster, C. Fransen, S. Freire, S. Fujimoto, H. Furuya, 512 M. Galbany-Casals, A. Gale, H. Galea, T. Gao, R. Garic, S. Garnett, S. Gaviria-Melo, S. Gerken, 513 D. Gibson, R. Gibson, J. Gil, A. Gittenberger, C. Glasby, H. Glenner, A. Glover, S.E. Gómez-514 Noguera, A.I. Gondim, B. Gonzalez, D. González-Solís, C. Goodwin, M. Gostel, M. Grabowski, C. Gravili, M. Grossi, J.M.. Guerra-García, J.M. Guerrero, R. Guidetti, M.D. Guiry, D. Gutierrez, 515 K.A. Hadfield, E. Hajdu, K. Halanych, J. Hallermann, B.W. Hayward, T.A. Hegna, G. Heiden, 516 E. Hendrycks, D. Hennen, D. Herbert, A. Herrera Bachiller, M. Hodda, J. Høeg, B. Hoeksema, 517 O. Holovachov, M.D. Hooge, J.N. Hooper, T. Horton, R. Houart, R. Huys, M. Hyžný, L.F.M. 518 Iniesta, T. Iseto, M. Iwataki, R. Janssen, D. Jaume, K. Jazdzewski, C.D. Jersabek, P. Jiménez-519 Mejías, P. Jóźwiak, A. Kabat, K. Kakui, Y. Kantor, I. Karanovic, B. Karapunar, B. Karthick, J. Kathirithamby, L. Katinas, N. Kilian, Y.H. Kim, R. King, P.M. Kirk, M. Klautau, J.P. Kociolek, 521 F. Köhler, K. Konowalik, A. Kotov, Z. Kovács, A. Kremenetskaia, R.M. Kristensen, A. Kroh, 522 M. Kulikovskiy, S. Kullander, E. Kupriyanova, A. Lamaro, G. Lambert, I. Laridon, D. Lazarus, 523 F. Le Coze, M. Le Roux, S. LeCroy, D. Leduc, E.J. Lefkowitz, R. Lemaitre, I.H. Lichter-Marck, 524 S.C. Lim, D. Lindsay, Y. Liu, B. Loeuille, A.N. Lörz, T. Ludwig, N. Lundholm, E. Macpherson, C. Mah, T. Mamos, R. Manconi, G. Mapstone, P.E. Marek, K. Markello, B. Marshall, D.J. Marshall, P. Martin, P. Martinez Arbizu, C. McFadden, S.J. McInnes, R. McKenzie, J. Means, J. Mees, H.H. Mejía-Madrid, K. Meland, K.L. Merrin, J. Miller, C. Mills, Ø. Moestrup, V. Mok-527 ievsky, T. Molodtsova, F. Monniot, R. Mooi, A.C. Morandini, R. Moreira da Rocha, C. Morrow, 528 J. Mortelmans, A. Müller, A.R. Muñoz Gallego, L. Musco, A.L.D.S. Nascimento, J.B. Nascimento, 529 G. Nesom, E. Neubert, B. Neuhaus, P. Ng, A.D. Nguyen, C. Nielsen, S. Nielsen, T. Nishikawa, 530 J. Norenburg, T. O'Hara, D. Opresko, M. Osawa, H.J. Osigus, Y. Ota, B. Páll-Gergely, J.L. Panero, D. Patterson, M. Pedram, P. Pelser, R. Peña Santiago, J.d.S.. Pereira, M. Perez-Losada, I. Petrescu, T. Pfingstl, W. Piasecki, D. Pica, B. Picton, J. Pignatti, J.F. Pilger, U. Pinheiro, A.B. Pisera, B. Poatskievick Pierezan, D. Polhemus, G.C. Poore, M. Potapova, R.A. Praxedes, V. Půža, G. Read, 534 M. Reich, J.D. Reimer, H. Reip, V. Resende Bueno, M. Reuscher, J.W. Reynolds, I. Richling, F. Rimet, P. Ríos, M. Rius, E. Rodríguez, D.C. Rogers, N. Roque, G. Rosenberg, K. Rützler, M. Saavedra, K. Sabbe, R. Sabroux, J. Saiz-Salinas, S. Sala, K. Samimi-Namin, S. Santagata, S. Santos, S.G. Santos, E. Sar, T. Saucède, L. Schärer, B. Schierwater, E. Schilling, A. Schmidt-Lebuhn, A. Schmidt-Rhaesa, S. Schneider, C. Schönberg, J. Schrével, P. Schuchert, C. Schweitzer, 538 J.C. Semple, A.R. Senna, A. Sennikov, C. Serejo, S. Shaik, S. Shamsi, J. Sharma, W.A. Shear, N. Shenkar, M. Short, J. Sicinski, D. Sidorov, P. Sierwald, D.K.F.d. Silva, E.S.S. Silva, E. Simmons,

540 541 542 543 544 545 546 546 547 548 549 550 551	F. Sinniger, C. Sinou, D. Sivell, H. Smit, N. Smit, N. Smol, M.V. Sørensen, J.F Souza-Filho, J. Spelda, W. Sterrer, H.M. Steyn, P. Stoev, S. Stöhr, E. Suárez-Morales, A. Susanna, C. Suttle, B.J. Swalla, S. Taiti, M. Tanaka, A.H. Tandberg, D. Tang, M. Tasker, J. Taylor, J. Taylor, K. Taylor, A. Tchesunov, E. Temereva, H. ten Hove, J.J. ter Poorten, K. Thirouin, J.D. Thomas, E.V. Thuesen, M. Thurston, B. Thuy, J.T. Timi, A. Todaro, J. Todd, X. Turon, P. Uetz, L. Urbatsch, J. Uribe-Palomino, E. Urtubey, S. Utevsky, J. Vacelet, D. Vachard, W. Vader, R. Väinölä, G. Valls Domedel, B. Van de Vijver, S.E. van der Meij, T. van Haaren, R.W. van Soest, A. Vanreusel, V. Venekey, T. Verhoeff, M. Vinarski, R. Vonk, C. Vos, A.A. Vouilloud, G. Walker-Smith, T.C. Walter, L. Watling, M. Wayland, T. Wesener, C.E. Wetzel, C. Whipps, K. White, U. Wieneke, D.M. Williams, G. Williams, R. Wilson, J. Witkowski, N. Wyatt, J. Xavier, K. Xu, J. Zanol, W. Zeidler, Z. Zhao, and A. Zullini. World register of marine species (worms). =https://www.marinespecies.org, 2024. URL https://www.marinespecies.org. Accessed: 2024-04-04.
552 553 554	Basit Alawode, Yuhang Guo, Mehnaz Ummar, Naoufel Werghi, Jorge Dias, Ajmal Mian, and Sajid Javed. Utb180: A high-quality benchmark for underwater tracking. In <i>Asian Conference on Computer Vision (ACCV)</i> , pp. 3326–3342, 2022.
555 556 557 558	Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In <i>Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization</i> , pp. 65–72, 2005.
559 560 561 562	Océane Boulais, Simegnew Yihunie Alaba, John E Ball, Matthew Campbell, Ahmed Tashfin Iftekhar, Robert Moorehead, James Primrose, Jack Prior, Farron Wallace, Henry Yu, et al. Seamapd21: A large-scale reef fish dataset for fine-grained categorization. In <i>Proceedings of the FGVC8: The</i> <i>Eight Workshop on Fine-Grained Visual Categorization, Online</i> , volume 25, pp. 2, 2021.
563 564 565	Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. <i>arXiv preprint arXiv:1906.07155</i> , 2019.
566 567	Jang Hyun Cho and Philipp Krähenbühl. Language-conditioned detection transformer. <i>arXiv preprint arXiv:2311.17902</i> , 2023.
568 569 570 571	Mark John Costello, Simon Claus, Stefanie Dekeyzer, Leen Vandepitte, Éamonn Ó Tuama, Dan Lear, and Harvey Tyler-Walters. Biological and ecological traits of marine species. <i>PeerJ</i> , 3:e1201, 2015.
572 573 574 575	Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
576 577 578 579	Roberto Danovaro, Joan Batista Company, Cinzia Corinaldesi, Gianfranco D'Onghia, Bella Galil, Cristina Gambi, Andrew J Gooday, Nikolaos Lampadariou, Gian Marco Luna, Caterina Morigi, et al. Deep-sea biodiversity in the mediterranean sea: the known, the unknown, and the unknowable. <i>PloS one</i> , 5(8):e11832, 2010.
580 581	Karan Desai, Gaurav Kaul, Zubin Aysola, and Justin Johnson. RedCaps: Web-curated image-text data created by the people, for the people. In <i>NeurIPS Datasets and Benchmarks</i> , 2021.
583 584 585 586	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. <i>arXiv preprint arXiv:2010.11929</i> , 2020.
587	Paul R Epstein, Rita R Colwell, and Timothy E Ford. Marine ecosystems. J. Onwhyn, 1993.
588 589 590	Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in 2021. <i>arXiv preprint arXiv:2107.08430</i> , 2021.
591 592 593	Danna Gurari, Qing Li, Abigale Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P. Bigham. Vizwiz grand challenge: Answering visual questions from blind people. 2018 <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 3608–3617, 2018. URL https://api.semanticscholar.org/CorpusID:3831582.

594 595 596	Benjamin S Halpern, Shaun Walbridge, Kimberly A Selkoe, Carrie V Kappel, Fiorenza Micheli, Caterina d'Agrosa, John F Bruno, Kenneth S Casey, Colin Ebert, Helen E Fox, et al. A global map of human impact on marine ecosystems. <i>science</i> , 319(5865):948–952, 2008.
597	
598	Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In <i>Proceedings of the</i>
599	<i>TEEE international conference on computer vision</i> , pp. 2961–2969, 2017.
600	Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi, Clipscore: A reference-
601	free evaluation metric for image captioning. arXiv preprint arXiv:2104.08718, 2021.
602	
604	JG Hiddink and R Ter Hofstede. Climate induced increases in species richness of marine fishes.
605	<i>Giobal change biology</i> , 14(3):453–460, 2008.
606	Jungseok Hong, Michael Fulton, and Junaed Sattar. Trashcan: A semantically-segmented dataset
607	towards visual detection of marine debris. arXiv preprint arXiv:2007.08097, 2020.
608	
609	Md Jahidul Islam, Chelsey Edge, Yuyang Xiao, Peigen Luo, Muntaqim Mehtaz, Christopher Morse,
610	sauman Sakib Enan, and Junacu Sallar. Semantic segmentation of underwater imagery: Dataset and henchmark. In <i>IEEE/RSI International Conference on Intelligent Robots and Systems (IPOS)</i>
611	pp. 1769–1776. IEEE, 2020.
612	
613	Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to
614	objects in photographs of natural scenes. In Proceedings of the 2014 conference on empirical
615	methods in natural language processing (EMNLP), pp. 787–798, 2014.
616	Faizan Faroog Khan Xiang Li Andrew I Temple and Mohamed Elhoseiny Fishnet: A large-
617	scale dataset and benchmark for fish recognition, detection, and functional trait prediction. In
618	Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 20496–
619	20506, October 2023a.
620	Evine Force When Ying Li Andrew I Townle and Mahamad Elhopping, Eisburg Allance colle
622	dataset and benchmark for fish recognition, detection, and functional trait prediction. In <i>IEEE/CVE</i>
623	International Conference on Computer Vision (ICCV), pp. 20496–20506, 2023b.
624	5 1 ()/11 /
625	Dahun Kim, Anelia Angelova, and Weicheng Kuo. Region-aware pretraining for open-vocabulary
626	object detection with vision transformers. In <i>IEEE/CVF Conference on Computer Vision and</i>
627	Pattern Recognition (CVPR), pp. 11144–11154, 2023.
628	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
629	Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In Proceedings
630	of the IEEE/CVF International Conference on Computer Vision, pp. 4015–4026, 2023.
631	Daniau Krishna, Vulca Zhu, Olivar Crath, Justin Jahnson, Kanii Hata, Jashua Kravitz, Stanhania
632	Chen Yannis Kalantidis Li-Jia Li David A Shamma et al Visual genome: Connecting language
633	and vision using crowdsourced dense image annotations. International journal of computer vision.
634	123:32–73, 2017.
635	
635	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-
638	Machine Learning (ICML) pp 12888–12000 PMLP 2022
630	machine Leanning (10mL), pp. 12000-12700. 1 MILA, 2022.
640	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
641	pre-training with frozen image encoders and large language models. International Conference on
642	Machine Learning (ICML), 2023a.
643	Junnan Li Dangyu Li Silvia Savaraga and Stavan Hai DI ID 2: hastatranning language image
644	nre-training with frozen image encoders and large language models. In <i>ICMI</i> , 2023b
645	pre duming with mozen image cheoders and large language models. In remil, 20230.
646	Lin Li, Eric Rigall, Junyu Dong, and Geng Chen. Mas3k: An open dataset for marine animal
647	segmentation. In <i>International Symposium on Benchmarking, Measuring and Optimization</i> , pp. 194–212. Springer, 2020.

648 Lin Li, Bo Dong, Eric Rigall, Tao Zhou, Junyu Dong, and Geng Chen. Marine animal segmentation. 649 IEEE Transactions on Circuits and Systems for Video Technology (TCSVT), 32(4):2303–2314, 650 2021. 651 Shijie Lian, Hua Li, Runmin Cong, Suqi Li, Wei Zhang, and Sam Kwong. Watermask: Instance 652 segmentation for underwater imagery. In Proceedings of the IEEE/CVF International Conference 653 on Computer Vision (ICCV), pp. 1305–1315, October 2023. 654 655 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In Text summarization branches out, pp. 74-81, 2004. 656 657 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 658 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer Vision-659 ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, 660 Part V 13, pp. 740–755. Springer, 2014. 661 Chongwei Liu, Haojie Li, Shuchang Wang, Ming Zhu, Dong Wang, Xin Fan, and Zhihui Wang. A 662 dataset and benchmark of underwater object detection for robot picking. In IEEE International 663 Conference on Multimedia and Expo Workshops (ICMEW), pp. 1–6. IEEE, 2021a. 664 665 Chongwei Liu, Zhihui Wang, Shijie Wang, Tao Tang, Yulong Tao, Caifei Yang, Haojie Li, Xing 666 Liu, and Xin Fan. A new dataset, poisson gan and aquanet for underwater object grabbing. IEEE 667 Transactions on Circuits and Systems for Video Technology (TCSVT), 32(5):2831–2844, 2021b. 668 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Neural 669 Information Processing Systems (Neurips), 2023a. 670 671 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances in neural information processing systems, 36, 2024. 672 673 Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei 674 Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for 675 open-set object detection. arXiv preprint arXiv:2303.05499, 2023b. 676 Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and 677 Alexander C Berg. Ssd: Single shot multibox detector. In European Conference Computer Vision, 678 pp. 21-37. Springer, 2016. 679 680 Xin Lu, Buyu Li, Yuxin Yue, Quanquan Li, and Junjie Yan. Grid r-cnn. In IEEE/CVF Conference on 681 Computer Vision and Pattern Recognition (CVPR), pp. 7363–7372, 2019. 682 Marta Miatta, Amanda E Bates, and Paul VR Snelgrove. Incorporating biological traits into conser-683 vation strategies. Annual Review of Marine Science, 13:421-443, 2021. 684 685 OpenAI. Introducing chatgpt. 2022. URL https://openai.com/blog/chatgpt. 686 OpenAI. Gpt-4 technical report, 2023. 687 688 Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg. Im2text: Describing images using 1 million 689 captioned photographs. In Neural Information Processing Systems (NIPS), 2011. 690 Rupert FG Ormond, JD Gagean, and Martin V Angel. Marine biodiversity. Patterns and Processes, 691 1997. 692 693 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic 694 evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pp. 311–318, 2002. 696 Allison L Perry, Paula J Low, Jim R Ellis, and John D Reynolds. Climate change and distribution 697 shifts in marine fishes. science, 308(5730):1912-1915, 2005. 698 Elvira S Poloczanska, Christopher J Brown, William J Sydeman, Wolfgang Kiessling, David S 699 Schoeman, Pippa J Moore, Keith Brander, John F Bruno, Lauren B Buckley, Michael T Burrows, 700 et al. Global imprint of climate change on marine life. Nature climate change, 3(10):919–925, 701 2013.

702 703 704 705	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International Conference on Machine Learning (ICML)</i> , pp. 8748–8763. PMLR, 2021.
706 707 708 709	Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Erix Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. <i>arXiv preprint arXiv:2311.03356</i> , 2023.
710 711 712 713 714 715	Cyrus Rashtchian, Peter Young, Micah Hodosh, and Julia Hockenmaier. Collecting image annotations using Amazon's Mechanical Turk. In Chris Callison-Burch and Mark Dredze (eds.), <i>Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk</i> , pp. 139–147, Los Angeles, June 2010. Association for Computational Linguistics. URL https://aclanthology.org/W10-0721.
716 717 718	Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 779–788, 2016.
719 720 721 722	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. <i>Advances in neural information processing systems</i> , 28, 2015.
723 724 725	Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 39(6):1137–1149, 2016.
726 727 728 729	Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. <i>arXiv preprint arXiv:2111.02114</i> , 2021.
730 731 732 733	Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 8430–8439, 2019.
734 735 736	Haozhan Shen, Tiancheng Zhao, Mingwei Zhu, and Jianwei Yin. Groundvlp: Harnessing zero-shot visual grounding from vision-language pre-training and open-vocabulary object detection, 2023.
737 738 739	Guolei Sun, Zhaochong An, Yun Liu, Ce Liu, Christos Sakaridis, Deng-Ping Fan, and Luc Van Gool. Indiscernible object counting in underwater scenes. In <i>IEEE/CVF International Conference on Computer Vision and Patern Recognition (CVPR)</i> , 2023.
740 741 742 743	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> , 2023.
744 745	Boyce Thorne-Miller. <i>The living ocean: understanding and protecting marine biodiversity</i> . Island Press, 1999.
746 747 748 749	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
750 751 752 753	Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 4566–4575, 2015.
754 755	Zhenyu Wang, Yali Li, Xi Chen, Ser-Nam Lim, Antonio Torralba, Hengshuang Zhao, and Shengjin Wang. Detecting everything in the open world: Towards universal object detection. In <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 11433–11443, 2023.

- Lewei Yao, Jianhua Han, Xiaodan Liang, Dan Xu, Wei Zhang, Zhenguo Li, and Hang Xu. Detclipv2: Scalable open-vocabulary object detection pre-training via word-region alignment. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 23497–23506, 2023.
- Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh International Conference on Learning Representations*, 2022.
- Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition* (*CVPR*), pp. 14393–14402, 2021.
- Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and
 Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023.
- Tiancheng Zhao, Tianqi Zhang, Mingwei Zhu, Haozhan Shen, Kyusong Lee, Xiaopeng Lu, and
 Jianwei Yin. Vl-checklist: Evaluating pre-trained vision-language models with objects, attributes
 and relations. *arXiv preprint arXiv:2207.00221*, 2022.
- Ziqiang Zheng, Jipeng Zhang, Tuan-Anh Vu, Shizhe Diao, Yue Him Wong Tim, and Sai-Kit Yeung. Marinegpt: Unlocking secrets of "ocean" to the public, 2023.
- Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li,
 Luowei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. Regionclip: Region-based language-image
 pretraining. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16793–16803, 2022.
- Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- Peiqin Zhuang, Yali Wang, and Yu Qiao. Wildfish: A large benchmark for fish recognition in the wild. In *ACM international conference on Multimedia (ACM MM)*, pp. 1301–1309, 2018.
- Peiqin Zhuang, Yali Wang, and Yu Qiao. Wildfish++: A comprehensive fish benchmark for multimedia research. *IEEE Transactions on Multimedia (TMM)*, 23:3603–3617, 2020.
- 792 793

A APPENDIX

796 A.1 DATA ANNOTATION

BBOX annotation. We develop a bounding box annotation platform as shown in Fig. 7. The left
side represents the image with the point prompts from the users. We embed the segment anything
model (SAM) as our labeling engine. The green dots indicate selected areas and the red dots indicate
unselected areas. The right slide is the labeled bounding box and corresponding instance mask output
automatically generated by SAM.

802 Captions refinement. We develop a user-friendly caption annotation platform as shown in Fig. 8. 803 On the left is the image and the object bounding boxes, and on the right is the intention description 804 of the object. We divided the images into 100 per subset and assigned them to experts so that the 805 duration of each consecutive work is not too long to ensure the quality of the annotation. We select one salient object to perform caption refinement and tag "Refined Positive" as shown in Fig. 8. The 806 generated captions are then distinguished into "Generated Positive" and "Generated Negative" for 807 the remaining objects. Fig. 9 is the example of "Generated Negative". We provide additional tags 808 from three aspects (Object, Relation, and Attribute) including 11 properties, Table 5 presents the 809 example and statistic.

	Properties	Example	Number
Object	Classification	This is a yellow <u>fish</u> . vs. This is a yellow <u>coral</u> .	6,875
	Background	The turtle is in the <u>ocean</u> . vs. The turtle is in the <u>sky</u> .	1,343
	Unexisting	The shark has a long tail. (there is no tail in the image)	3,264
Relation	Spatial	This fish is <u>under</u> the coral. vs . This fish is <u>on</u> the coral.	816
	Action	The penguin is <u>walking</u> . vs . The penguin is <u>sitting</u> .	938
Attribute	Size	The shark is large. vs. The shark is <u>small</u> .	271
	Color	This is a <u>yellow</u> fish. vs. This is a <u>blue</u> fish.	2,031
	Shape	This is a <u>oval</u> seashell. vs. This is a triangle seashell.	312
	Texture	The seashell is <u>smooth</u> . vs. The seashell is <u>rough</u> .	321
	Material	The fish is probably made of <u>plastic</u> .	316
	Counting	There are <u>three</u> penguins vs. There are four penguins.	831

Table 5: The detailed explanations of the constructed 11 attributes and corresponding data statistics for the generated negative captions.

A.2 EXAMPLES AND DATA DIVERSITY

Examples. We provide some image examples with the detailed bounding box instance caption annotations in Fig. 10. We encourage the readers to pay more attention to the generated instance captions. The instance captions describe the appearance of the object instance, action, event, the relationship between the selected instance with other instances, and more advanced biological traits.

Diversity and data composition. We provide the illustration about the data diversity of our constructed MarineMaid dataset in Fig. 11. We only provide some images from some categories for better illustration.



Figure 7: Screenshot of the BBOX annotation platform. *Left:* input point prompt. *Right:* the labeled instance BBOX and mask.

MORE EXPERIMENTS В

B.1 DATA SPLIT

We implement a consistent splitting strategy for each dataset: class-level, intra-class, and inter-class. For the training set, 80% of the images containing objects from seen categories are randomly sampled. The remaining 20% of seen objects, along with all unseen objects, are allocated to the validation set. To assess performance on both seen and unseen objects, the validation set is further divided into val_seen and val_unseen based on categories. Images containing both seen and unseen objects are manually reassigned to the validation set, resulting in duplicated images in both val_seen and val_unseen, each with different annotations. Statistics and details can be found in Table 6.

885

886 887



Figure 8: Screenshot of our developed caption refinement platform for generating the "*Refined Positive*". The domain experts are required to modify and edit the accurate and detailed biological traits for the selected marine instance.

Text Label



Figure 9: Screenshot of our developed caption refinement platform for generating the "*Refined Negative*". The annotators are asked to provide additional attribute annotations (wrong types) for the negative captions.

Table 6: Data split for performing the open-vocabulary object detection. We provide detailed data
 splitting under each setting.

913													
			Class-	-Level			Intra	-Class			Inter	-Class	
914		Train	vəl	val	val	Train	val	val	val	Train	vəl	val	val
915		ITain	vai	(seen)	(unseen)	IIaiii	vai	(seen)	(unseen)	IIaiii	vai	(seen)	(unseen)
016	# of Images	9,291	5,298	2,743	2,963	9,312	5,301	2,746	2,963	8,163	6,062	2,943	4,023
910	# of BBOX	27,560	14,540	9,071	5,469	28,166	13,992	8,523	5,469	23,561	17,920	9,960	7,960
917	# of Captions	22,739	11,897	6,699	5,198	22,709	11,984	6,786	5,198	19,232	14,790	7,527	7,263

This is an image of a fish in someone's hand. The fish appears to have a long, slender body with a flat head. Its gills are visible. The fish has large and black pectoral fins.

This image shows a damselfish with yellow eyes. The fish has light colored face but darker body and large fins. It is surrounded by coral and rocks in the background.

This figure appears to be a photograph of a fish with black and white stripes and white spots. The fish is swimming in an aquarium or other body of water.

an he was



918



923

924

- 925
- 926
- 927 928
- 929
- 930

931

932

933



935 936

937





942

943



945

946

947 948

949



952

953

954

955



957 958



960

961 962

963



966

967

upwards

This image depicts a small red squat lobster with large eyes looking out from inside of a rock or coral in the ocean. It is red in colour with fuzzy limbs.



The object in this figure is a shark. It appears to be lying on its bottom and the sand. The shark's body is brown and it has 2 barbels in front of its mouth.



This is an image of a seashell. The shell has a brown and white pattern on it and appears to be sitting on top of a sandy beach.

The image shows a beluga whale with its head above



The image shows a dead frilled shark with its mouth open, exposing its rows of sharp teeth which run perpendicular to its jaw. The shark has a large green eye and brown, slimy textured skin.



This image appears to show a close-up view of a colorful scorpionfish with bright orange eyes and pink and orange coloration on its body. It appears to be underwater with corals in the background.



This is an image of an orange fish with a white stripe down its back with half its body protruding from a sea anemone. The anemone has greenish-yellow tentacles all facing the same direction as the fish.



The image shows a close-up view of a coconut crab that is facing towards the camera. It has giant, muscular and spiny claws that are brown in color and open. There are long, outstretched walking legs in blueish color.



This is a close-up view of a coelacanth swimming between rocks. It has a large body in blueish green with irregular white spots and lobed pectoral fins that extend away from the body.



This image appears to be an emperor angelfish, which is a species of marine fish that is found in the pacific and indian oceans. The fish has a blue and yellow body with stripes and a long, yellow tail.



The image shows a close-up view of a bobtail squid on a black background. It has a pair of large eyes, and a body and short tentacles in green colors with many brown spots.



The image is of a dolphin swimming on the surface of the ocean. The dolphin has dark gray pattern on its body. It is located at the left of this image. It is swimming with another dolphin.



The object in the illustration is a frilled shark. It has a long narrow body, visible gills behind its head and a green eye. Its teeth are prominent and sharp. The illustration is on a white background.



This is an image of a small triplefin fish with a large head and small, yellow body. The fish appears to be hovering over rocks that are covered in light-coloured alread algae



This is the top view of a basket star on the left side of the image. It has long, curly arms with many branches. The arms appear to be entangled with each other.



This image shows a brightly colored fish with an orange body, pink fins, and a blue edge to it's fins and tail. It is likely to be a species of serranidae. It is swimming in front of a dark background.

Figure 10: The example images with the bounding box annotations and instance captions from our MarineMaid dataset. Best viewed in color.

968 969

970



Figure 11: We present the data distribution of our MarineMaid dataset at the "Class" level. We also provide images from some selected Classes for illustration.

998 999 1000

Table 7: Results of MiniGPT-4 under two settings on our MarineMaid dataset.

Method	CLIPScore↑	RefCLIPScore↑	CIDEr↑	BLUE-4↑	METEOR↑	Rouge↑
Vanilla	74.48	73.43	5.72	7.18	16.90	28.03
Fine-tuned	77.96	77.51	17.36	14.79	16.90	33.71

1001 B.2 INSTANCE CAPTIONING

Due to the constraint of the computational power, we select the representative MiniGPT-4 (with LLaMA2 7B version) to do the fine-tuning on our MarineMaid dataset. Please note that the testing set is withheld for evaluation purposes. We finetuned the MiniGPT-4 on our dataset on 4 NVIDIA A100-40GB for 5 epochs and we set other training parameters to follow the same as its original paper. We report the experimental results under the two settings (vanilla and fine-tuned) in Table 7. We observe that further fine-tuning could help improve the instance understanding performance. But there is still large room for further improvement. The sole fine-tuning cannot fully solve our problem and domain-specific design and modifications are required.

1011 1012 B.3 GROUNDING

1013 Following a similar experimental setting, we select GroundindDINO to perform the grounding 1014 experiments. To guarantee the proper configuration of the evaluation environment settings, we 1015 meticulously adhered to the instructions and evaluation program provided by the original official 1016 implementations. We further fine-tune GroundingDINO on our training dataset to measure the 1017 performance. We use the same pre-trained model (MM-GDINO-B*) to optimize the model on our dataset. The results are reported in Table 8. We observe an observable improvement in R@1 across 1018 all validation settings, though there is a decline in performance in R@5 and R@10. This indicates 1019 that after fine-tuning, the model becomes more proficient at identifying the target object within the 1020 image based on the query but also detects additional regions that do not correspond to the ground 1021 truth bounding box. However, we also acknowledge the performance drop of the R@10. We attribute 1022 such performance drop to the specific feature of our constructed MarineMaid dataset (we only aim to 1023 ground one instance based on the captions). 1024

1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051

Table 8: Performance comparison of GroundingDINO under two settings: Vanilla and Fine-tuned.

		Seen			Unseen		
		Class-Level	Inter-Class	Inter-Class	Class-Level	Inter-Class	Inter-Class
Vanilla	R@1	38.8	37.5	37.2	46.8	46.8	45.6
	R@5	67.3	65.5	66.2	76.9	76.8	74.5
	R@10	78.0	76.0	76.8	84.9	84.9	83.3
Fine-tuned	R@1	65.7	62.1	64.5	71.6	70.6	70.3
	R@5	71.9	66.1	63.5	76.0	73.6	71.4
	R@10	74.2	67.7	64.1	78.5	74.8	71.7