

The Coralscapes Dataset: Semantic Scene Understanding in Coral Reefs

Jonathan Sauder^{*,1,2}, Viktor Domazetoski^{*,3,4}, Guilhem Banc-Prandi²,
Gabriela Perna², Anders Meibom^{2,5}, Devis Tuia¹

¹Environmental Computational Science and Earth Observation Laboratory, École Polytechnique Fédérale de Lausanne, Switzerland

²Laboratory for Biological Geochemistry, École Polytechnique Fédérale de Lausanne, Switzerland

³Centre for Ecology and Conservation, University of Exeter, United Kingdom

⁴School of the Environment, The University of Queensland, Australia

⁵Center for Advanced Surface Analysis, University of Lausanne, Switzerland



Abstract

Coral reefs are declining worldwide due to climate change and local stressors. To inform effective conservation or restoration, monitoring at the highest possible spatial and temporal resolution is necessary. Conventional coral reef surveying methods are limited in scalability due to their reliance on expert labor time, motivating the use of computer vision tools to automate the identification and abundance estimation of live corals from images. However, the design and evaluation of such tools has been impeded by the lack of large high quality datasets. We release the Coralscapes dataset, the first general-purpose dense semantic segmentation dataset for coral reefs, covering 2075 images, 39 benthic classes, and 174k segmentation masks annotated by experts. Coralscapes has a similar scope and the same structure as the widely used Cityscapes dataset for urban scene segmentation, allowing benchmarking of semantic segmentation models in a new challenging domain which requires expert knowledge to annotate. We benchmark a wide range of semantic segmentation models, and find that transfer learning from Coralscapes to existing smaller datasets consistently leads to state-of-the-art performance. Coralscapes will catalyze research on efficient, scalable, and standardized coral reef surveying methods based on computer vision, and holds the potential to streamline the development of underwater ecological robotics.

I. INTRODUCTION

Coral reefs are under unprecedented stress from both global anthropogenic climate change and local human activities such as overtourism, pollution, or destructive fishing practices [53]. Coral reefs, which host more than a third of all marine biodiversity on less than 0.1% of the world’s oceans’ surface [30], are among the most vulnerable ecosystems on the planet: under the current greenhouse gas emission trajectory, more than 99% of warm-water coral reefs are projected to lose significant area or suffer local extinction [54], [26]. This could cascade into catastrophic impacts for the more than 500 million people that rely on coral reef-related ecosystem services through food security, coastal protection, or tourism [58].

However, there are regions, species, and genotypes within species that are more heat-resistant, holding promise to withstand the projected ocean warming induced by climate change [13]. The reefs of the northern Red Sea, for example, are projected to resist up to 5°C of warming [29], [77], [45], [60], [71]. These heat-resistant reefs must be protected from local stressors to ensure the survival of coral reef ecosystems on the planet. To inform effective conservation or restoration strategies, reefs must be monitored at the highest spatial and temporal resolution. As remote sensing methods are limited in resolution, constrained to very shallow areas, and rely on good atmospheric,

* Equal contribution

weather, and water conditions [49], [6], [18], [41], *in situ* surveying remains the de facto method for monitoring coral reefs.

Computer vision emerges as a promising tool to scale coral reef monitoring, through tackling the bottleneck imposed by the needed expert time for identifying corals and estimating their abundance on images, as well as the development of visual simultaneous localization and mapping (SLAM) systems that can be used on autonomous underwater vehicles (AUVs) and other robots. However, existing datasets in the coral reef domain are limited to narrow domains such as orthomosaics or photo quadrats [5], [27], to sparse point labels [8], [9] or image classification [68], [39], or in their size and diversity, inhibiting the widespread deployment of computer vision in coral reef science and conservation.

The main challenge to creating large-scale classification/segmentation datasets in coral reefs is the annotation process. Annotating coral reefs requires trained experts and cannot be done by crowd-sourced click-workers. Corals display high morphologic plasticity [5], and the appearance of corals can vary drastically between biogeographic regions. In general, precise phylogenetic taxonomy of coral species or even genera can simply not be determined from images, even in close-up scenarios in good conditions, limiting the classification hierarchy to morphological attributes, such as growth forms (e.g. branching, meandering, massive). To add to the challenge, corals may appear bleached or dead, followed by various stages of degradation (e.g. overgrown by algae), or decomposing into rocky substrate or rubble. Besides corals, reefs host a wide range of life, including fish, invertebrates (sea urchins, sea cucumbers, anemones, sponges etc.), plants (algae, seagrass), and many other organisms. Finally, the identifiability of benthic classes in reefs is strongly impacted by the degradation of color and induced blur from the water column. Tackling these challenges to create a consistently labeled large-scale dataset requires a concerted effort by trained expert annotators.

We propose the Coralscapes dataset for semantic scene understanding in coral reefs. Coralscapes serves two main purposes: i) to drive the development of computer vision applications in reef conservation, including automated surveying methods and underwater robotics, and ii) to allow benchmarking of state-of-the-art semantic segmentation on a domain that is underrepresented in large pre-training datasets.

The Coralscapes dataset is:

- The first general-purpose semantic segmentation dataset in the coral reef domain, including a wide range of reefs from thriving to heavily bleached/dead, as well as diverse camera angles and distances to the substrate. Unlike previous datasets, Coralscapes is not restricted to close-ups of corals, orthomosaics, or photo quadrats.
- The largest expert-annotated semantic segmentation dataset spanning 2075 images at 1024×2048 px resolution, containing 174k polygons over 39 classes, labeled in a consistent and speculation-free manner.
- The first dataset in the coral reef domain designed for fair evaluation of machine learning approaches, enforcing a spatial train/test split on images gathered from 35 dive sites in 5 countries in the Red Sea.

II. RELATED WORK

A. Large Datasets for Semantic Segmentation

Large benchmark datasets have played a key role in deep learning's rapid advancement in semantic segmentation many other computer vision tasks. Early datasets such as PASCAL VOC [28] provided pixel-wise annotations for object categories, even before deep learning was ubiquitous. The introduction of Cityscapes [23], a dataset of 5000 images at 1024×2048 px with fine-grained semantic annotations in urban driving scenarios, set a precedent for large-scale semantic segmentation datasets. Notably, Cityscapes became the standard benchmark for evaluating segmentation networks, influencing the design of widely used neural network architectures such as the DeepLab [19] family of models [20], [21] and Pyramid Scene Parsing Networks [83].

Following the success of Cityscapes, several semantic segmentation datasets have been developed: ADE20K [85] includes a broad range of everyday environments, including street scenes, on many more labeled images than Cityscapes. Similarly, Mapillary Vistas [57] extended both the scale of urban scene segmentation and its scope by introducing images captured under diverse lighting, weather, and geographical conditions. Datasets like COCO-Stuff [17] and LVIS [33] explore object segmentation in a more general variety of contexts.

A multitude of semantic segmentation datasets in domains that commonly require experts have been proposed, such as the BRATS [55], LITS [14] and ACDC [11] datasets in the medical imaging domain, and the DeepGlobe [25], Potsdam [1], and Vaihingen [2] datasets in remote sensing. Expert-domain datasets are significantly smaller than general-purpose semantic segmentation datasets due to the limited availability of expert time for annotation.

Dataset	# Images	Annotation Type	# Annotations	# Annotated Pixels	# Classes	Comment	Spatial Split	Openly Available
RSMAS [73]	776	Image	776	-	14	-	No	Yes
Jamil et al. [39]	1582	Image	1582	-	3	-	No	Yes
Raphael et al. [68]	5000	Image	5000	-	11	Coral only	No	No
DeOhloNosCorais [31]	1411	Image + Foreground Map	1411	7.08M	21	-	No	Yes
Eilat [10]	212	Sparse	1123	-	8	-	No	Yes
King et al. [43]	1807	Sparse	9511	-	10	-	No	No
Moorea Labeled Corals [8]	2055	Sparse	400,000	-	9	Only photo quadrats	No	Yes
Pacific Labeled Corals [9]	5090	Sparse	251,988	-	20	Only photo quadrats	No	Yes
Benthos-2015 [12]	9874	Sparse	407,968	-	148	Ortho-photos	Yes	Yes
Seaview [32]	11387	Sparse	859,870	-	228	Ortho-photos	Yes	Yes
BenthicNet [51] *	29603	Sparse	1,390,262	-	262	*Subset of coral sites	Yes	Yes
Benthos [82]	4	Semantic Segmentation	4500	Unknown	8	Orthomosaics	No	Yes
UCSD Mosaics [27]	16	Semantic Segmentation	54055	1.29B	35	Orthomosaics	No	No
CoralSCOP [84]	41297	Segmentation Masks [†]	330,144	17.56B	136	[†] Coral masks only	No	Request
Coralscapes (Ours)	2075	Semantic Segmentation	174,077	3.36B	39	General purpose	Yes	Yes

Table I: Comparison of existing datasets for computer vision in coral reefs in terms of annotation type, size, classes, availability of disjoint geographic locations allowing a spatially disjoint train/test split, and open availability.

B. Coral Reef Classification & Segmentation

In the context of coral reefs, existing expert-labeled datasets provide either only image-level labels [73], [39], [68], [31] on close-up images of corals, a few sparse labeled pixels on entire images [10], [43], [8], [9], or focus on segmentation in restricted domains such as photo quadrats or orthomosaics [82], [27]. Table I provides a comparison of existing datasets, including the BenthicNet [51] dataset (subselected for coral images, i.e. all sites with at least one hard coral label and shallower than 60m), and CoralSCOP [84], which provides a model similar to Segment Anything [44] trained on annotated masks of corals (but no other benthic classes).

Most existing datasets do not split the training, validation, and testing data according to different geographic sites, leading to a likely overestimation of the classification performance compared to real-world applications in unseen reefs. Furthermore, the larger, but sparsely annotated datasets [51], [12], [32] are essentially crowd-sourced from various annotation teams, leading to a large and noisy label set, that is challenging to transfer to new biogeographic regions. Lastly, except for the SUIM dataset [38], which provides dense segmentation labels focusing on robot-human interaction on diving missions and has no dedicated classes for live corals, no existing segmentation datasets are *general-purpose* in the sense of diverse camera angles, distances to the substrate, environmental conditions (water turbidity, color, lighting), and with all relevant classes visible on the images annotated.

The absence of a high-quality general purpose segmentation dataset means that widely used machine learning applications in coral reefs are restricted to point- or patch-wise classification [8], [22], [56], [3], [35], and raised interest in automatic expansion of sparse labels into dense masks [5], [82], [64]. The Coralscapes dataset aims to fill this gap, providing a dataset built for dense segmentation in general purpose reef scenarios.

III. DATASET

A. Data Collection

All imagery was collected during scuba dives at 35 sites in coral reefs of Djibouti, Eritrea, Sudan, Jordan, and Israel using GoPro Hero 10 cameras. The annotated images are video frames originally taken at 1080×1920 px resolution and 30 frames per second in the linear lens setting. The videos were taken in a diverse set of reefs: from thriving reefs with very high live coral cover to devastated environments with severe bleaching and coral mortality, from high visibility in good lighting conditions to turbid scenes in low-light settings. The videos were acquired during transect surveying campaigns of the Transnational Red Sea Center hosted at the École Polytechnique Fédérale de Lausanne (EPFL). Therefore, Coralscapes includes images including transect tools, transect lines, divers, and some human-made structures commonly encountered in reef monitoring settings. To protect the coral reefs in Coralscapes from overtourism or illegal poaching activities, and in accordance with local research permit authorities, the location of the sites is withheld and instead replaced by an ID.

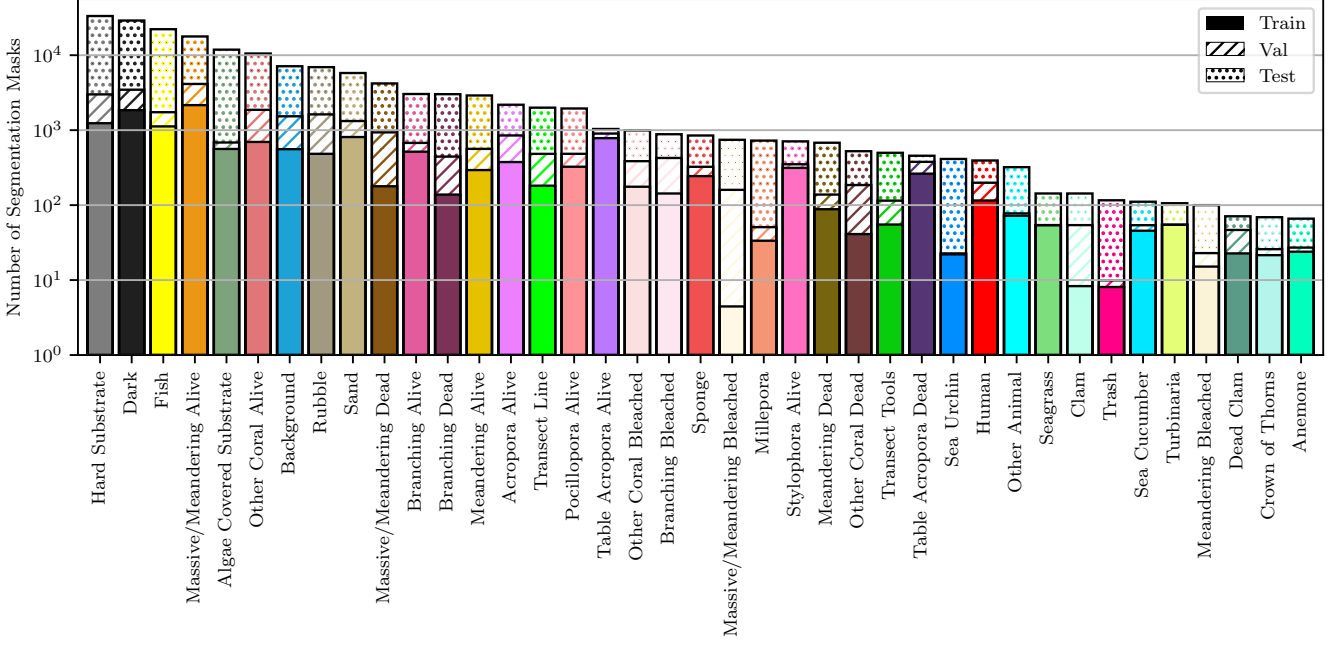


Fig. 1: Number of annotated segmentation masks per class in the Coralscapes dataset splits for each of the 39 classes (shown with linear proportions on logarithmic scale).

B. Annotation

The main challenge in creating a high-quality dataset for coral reef segmentation lies in the complexity of the annotations. Besides the visual degradation induced by the water column, many benthic classes are difficult to differentiate even in good conditions and by experts: classifying corals to a high taxonomic level from images is complex because of their strong morphological plasticity [76], [5]. Although some genera can be identified by visually distinct features, many genera can simply not be discerned from casual imagery, even by domain experts familiar with the biogeographic area. In such cases, the growth form of the coral (e.g. branching, massive) is commonly reported, as this determines the main ecological function of the coral within the reef ecosystem. Live corals are not the only source of ambiguity: fine algal mats only become apparent when viewed closely and are indistinguishable from non-algae-covered substrates at a larger distance. Additionally, when hard corals die, their calcium carbonate skeletons start degrading over time, and transform into hard substrate that may become itself covered in algae or break down into rubble. As a result, these substrates often fit multiple label classes depending on their stage of degradation. In spite of the inherent ambiguities, these classes remain of real interest to coral surveyors that monitor how a reef’s state changes with time, and should be identified by an automatic tool. At the same time, there should be no *speculative* assignment of polygons to classes that may be incorrect.

To align to all those requirements, the annotation was done in a speculation-free manner, enforcing a conservative semantic segmentation in terms of taxonomic depth. For example, when a branching coral was too far away from the camera to decide with certainty whether it is alive or dead, it was labeled as ‘background’. Viewed closer in another video frame, the same coral may be assigned to its growth form (branching, massive, etc.) and whether it is alive or dead. A precise genus label was assigned only when the genus was clearly distinguishable. In-depth details of the annotated classes is given in Appendix B.

A set of 39 classes was chosen to cover the most prominent and most important benthic classes present in the images (Figure 1). This label set includes 10 visually distinct live coral classes, 5 dead coral classes, 4 bleached coral classes, 6 non-coral invertebrates, 4 substrate classes, and classes ‘human’, ‘transect line’, ‘transect tools’, ‘fish’, and ‘trash’. The ‘dark’ and ‘background’ classes are assigned to pixels that can not be classified due to either poor lighting or being too far from the camera (and thus blurred or colorless). While the exact delineation of ‘dark’ and ‘background’ is subjective, their inclusion is necessary to reduce speculative labels.

Annotation of the selected 39 classes with segmentation masks was done by four experts on coral reefs who were

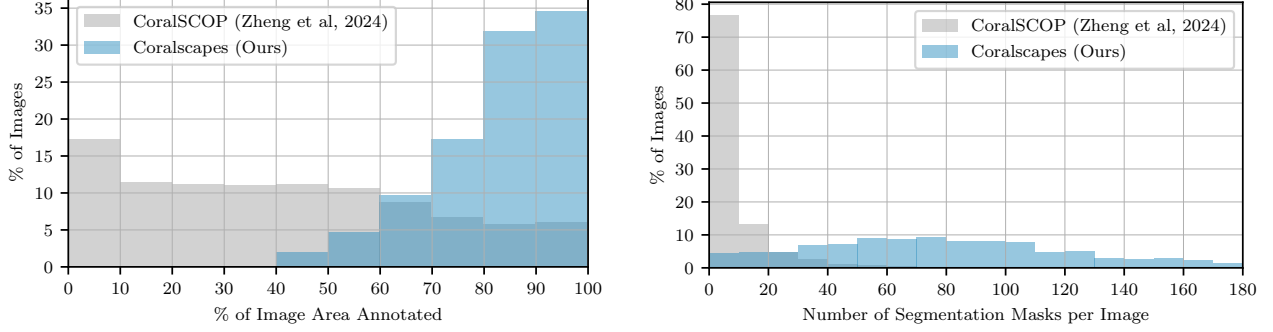


Fig. 2: Histogram of the image area covered by annotations (left) and of the number of polygons present per image (right), highlighting the complexity of semantic segmentation in Coralscapes, compared to the CoralSCOP dataset [84].

also present during data collection, and six additional trained annotators. The specific image frames for annotation were selected from over 200 hours of video material in order to capture the diversity of the scenes, while obtaining a good coverage of the 39 classes. The CVAT annotation interface [24] was used, which allowed annotators to interact with the Segment Anything Model [44]. While this helped to segment some objects with salient boundaries, it proved to be largely unhelpful for classes with less apparent boundaries and those appearing as small polygons, where polygons were drawn by hand. The annotation process emphasized high-quality annotations, meaning that there was a focus on annotating each frame with many polygons and on covering much of the image area, instead of annotating many images sparsely, as shown in Fig. 2. All polygons were visually verified to be correctly delineated and assigned the correct class.

In terms of class presence, the annotations of Coralscapes capture the strong imbalance in reef scenes: despite making an effort to prioritize the selection of frames with rare classes, the most common class has over 400 times more annotated polygons than the rarest class, as shown in Figure 1. Similarly, the granularity and size of the visible benthic classes varies dramatically: even though a class may appear very often, its total number of annotated pixels may still be low because it often appears in small polygons, and vice versa, as highlighted in Figure 3. For example, the median size of a ‘fish’ polygon is more than 120 times smaller than the median ‘background’ polygon. This imbalance in terms of representation and scale contributes to making Coralscapes a challenging benchmark for semantic segmentation. Additional statistics of Coralscapes (pixel counts, class frequency at image level) are provided in Appendix A.

C. Dataset Structure

We split the dataset spatially by reef site to allow a fair evaluation, resulting in a training set of 1517 images (27 sites), a validation set of 166 images (3 sites), and a test set of 392 images (5 sites). These splits were selected so that the classes are well represented in all three splits, which is shown in the proportions in Figure 1 (‘Turbinaria’ and ‘Seagrass’ are only present in two sites each, so they are omitted from the validation set).

To facilitate the usability of Coralscapes, the structure of the dataset mimics that of the widely used Cityscapes benchmark: images are provided as 8-bit PNG images resized to 1024×2048px resolution, and the 19 preceding and 10 trailing video frames are available for download (at 30 FPS). We make Coralscapes easily available: frames and their segmentation masks are available on Huggingface*, from where they can be loaded to Python in one line, and on Zenodo†, where we also provide the preceding and trailing frames, the frames at original 1080×1920px resolution, and benchmarked model checkpoints.

*<https://huggingface.co/datasets/EPFL-ECEO/coralscapes>

†<https://zenodo.org/records/15061505>

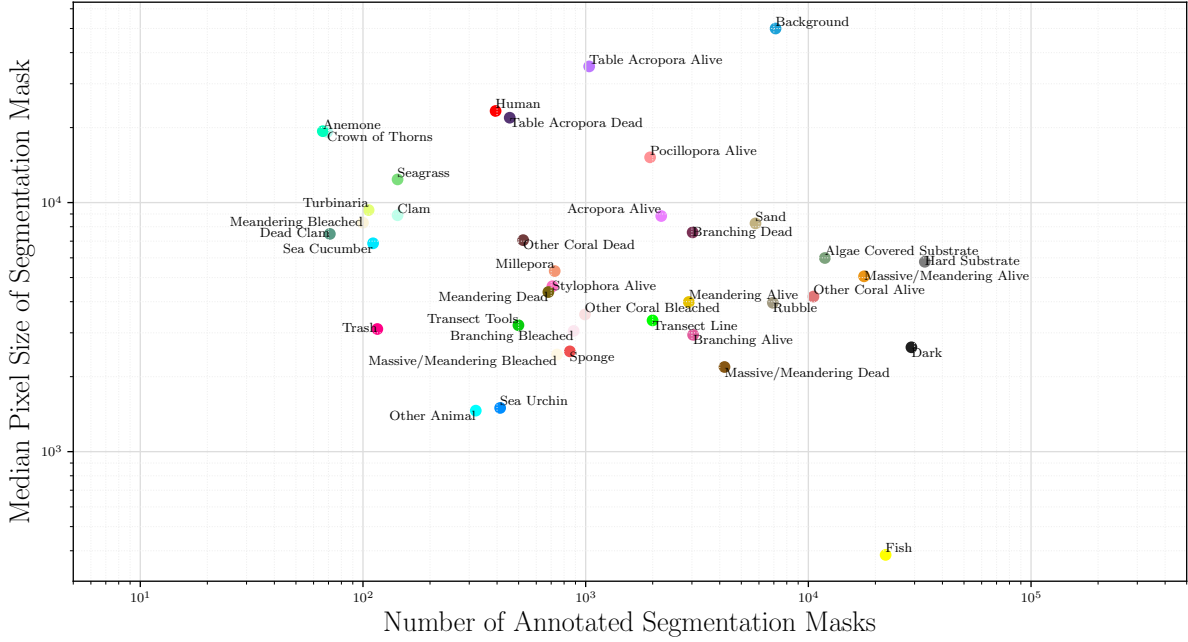


Fig. 3: Median size (in pixels) of an annotated polygon for the classes of Coralscapes, plotted against the number of annotated polygons. This highlights the challenge of segmenting classes that require the global image structure to segment correctly, as well as small fine-grained classes in the same dataset.

IV. BENCHMARKING

We benchmark a diverse set of commonly used semantic segmentation architectures to establish baseline performances for Coralscapes. Our evaluation includes both convolution-based and transformer-based models, covering a wide range of network complexities and training techniques. For convolutional architectures, we evaluate UNet++[86] and DeepLabV3+[21], both using a ResNet-50 backbone. For transformer-based architectures, we benchmark SegFormer [81] with both MiT-B2 and MiT-B5 backbones, each trained with and without Low-Rank Adaptation (LoRA) [36]. Similarly, we use DPT [66] with DINOv2 [59] backbones at two encoder sizes (Base and Giant), again with and without LoRA fine-tuning. To further explore the capabilities of self-supervised vision transformers, we also include a linear segmentation head trained on top of DINOv2-Base features. The ResNet and MiT backbones are initialized with ImageNet-pretrained weights, while the DINOv2 backbones are pretrained on a large collection of images in a self-supervised fashion.

Whenever available, we follow the original authors’ training procedures on Cityscapes in terms of hyperparameters, augmentations, and inference/evaluation strategy, unless otherwise specified. Implementation details for all

Method	Test Accuracy	Test mIoU
UNet++ - ResNet50	75.593	42.906
DeepLabV3+ - ResNet50	78.171	45.720
SegFormer - MiT-B2	80.904	54.682
SegFormer - MiT-B2 (LoRA)	80.987	51.911
SegFormer - MiT-B5	80.939	<u>55.031</u>
SegFormer - MiT-B5 (LoRA)	80.863	52.775
Linear - DINOv2-Base	81.339	52.478
DPT - DINOv2-Base	78.124	44.203
DPT - DINOv2-Base (LoRA)	80.053	47.233
DPT - DINOv2-Giant	80.691	50.643
DPT - DINOv2-Giant (LoRA)	<u>81.701</u>	54.531
SegFormer - MiT-B5 (Stronger Augmentations)	82.761	57.800

Table II: Quantitative results on Coralscapes segmentation, best shown in **bold**, second best underlined.

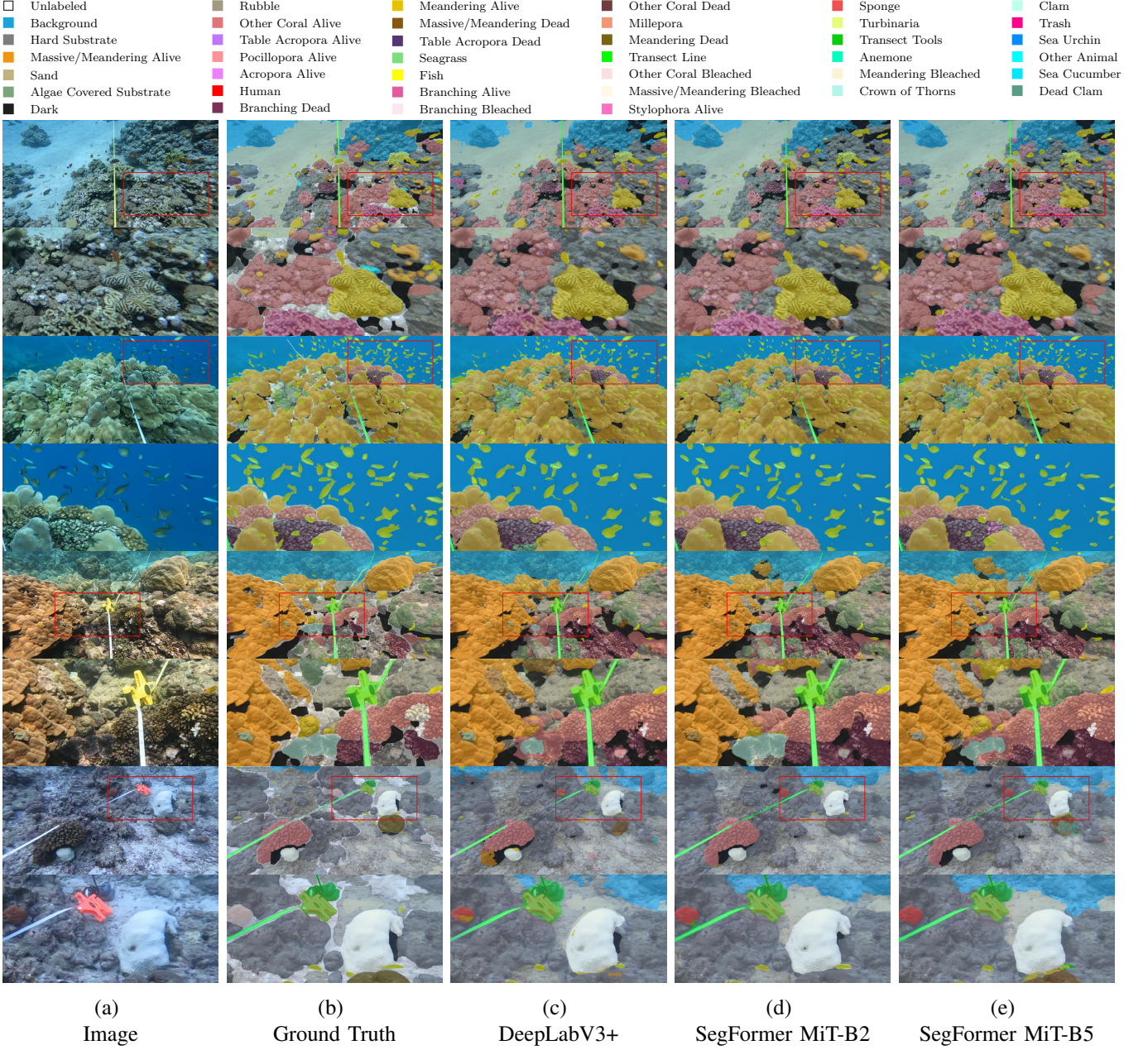


Fig. 4: Qualitative samples from the Coralscapes test set. Additional samples provided in the Appendix.

benchmarked models are provided in Appendix B. To choose the number of training epochs, we train each model on the training set only and evaluate on the validation set. We then re-train each model on the train+validation sets with the number of epochs leading to the highest validation mIoU (smoothed across 10 evaluation epochs in order to reduce the effect of noise). Quantitative results are shown in Table II, where we report performance of the models averaged over last three evaluated epochs, in order to reduce noise without training an ensemble of classifiers or cherry-picking results. A qualitative comparison of models is shown in Fig 4.

While the SegFormer models showed the best performance in terms of mIoU, the models with a DINOv2 backbone showed the best performance in terms of pixel accuracy. This likely stems from the fact that the DINOv2 models predict the segmentation maps for the image at once, where the global structure of the image helps determine the slightly subjective delineation of ‘background’, as opposed of the strided patch predictions from SegFormer, which processes the image at higher resolution and can predict the fine-grained classes better. We find that deviating from the training procedure of SegFormer on the Cityscapes dataset by increasing the strength of data augmentations, the performance can be further improved, yielding the best model.

V. TRANSFER LEARNING & APPLICATIONS

In this section, we evaluate the potential of transferring models trained on Coralscapes to three downstream tasks.

UCSD Mosaics Transfer Learning We evaluate models trained on Coralscapes in a transfer learning setting on the UCSD Mosaics dataset [27], which consists of 16 orthomosaics densely annotated with 33 benthic classes and one class for the remaining substrate. We follow [5], [64], [65] in dividing the orthomosaics into patches of size 512×512 px, and remove patches with corrupted ground-truth masks [64], leaving 3974 training and 696 test patches. We evaluate both in a dense segmentation setting, as well as in a commonly performed sparse-to-dense setting, in which the training annotations are sparsified, but evaluation is performed on the unchanged dense labels. In existing works [5], [64], [65], this setting is tackled by a two-step procedure, where in a first step, algorithms propagate the sparse labels into denser masks, and in a second step, a segmentation neural network is trained on the obtained denser masks. Our approach is conceptually simpler: we train segmentation models directly on the sparsified labels, ignoring the unlabeled pixels. We select 5, 10, 25, and 300 pixel locations per orthomosaic patch uniformly at random, comparing models pre-trained on Coralscapes (with the last layer re-initialized) against off-the-shelf models with a pre-trained backbone. The results, shown in Table III, show that pre-training on Coralscapes substantially improves segmentation with DeepLabV3+, particularly in highly sparse label regimes (5, 10, 25). When fine-tuning all parameters of a SegFormer MiT-B2, pre-training on Coralscapes improves performance in the sparse regimes, but not in the less sparse or dense (“All” in Table III) settings. When training only the lightweight decoder on the MiT-B5 encoder weights obtained by pre-training on Coralscapes, we observe a consistent performance increase compared to an ImageNet pre-trained encoder. Interestingly, even an off-the-shelf DeepLabV3+ trained on sparse randomly located labels outperforms existing two-step procedures in sparse label regimes, despite [65] showing that selecting pixels on a grid or with a human in the loop instead of randomly is beneficial.

Crown of Thorns Starfish Survey *Acanthaster planci*, colloquially known as the Crown of Thorns starfish (COTS), is a common reef inhabitant in the entire Indo-Pacific, where it preys on hard corals. While in normal population sizes, COTS are a healthy part of the ecosystem, there are outbreaks of COTS, in which their population increases dramatically and the hard coral cover declines, which can lead to devastating consequences for entire reefs [34], [63]. Interventions against COTS outbreaks are most effective in the early stages of an outbreak [62], [79]: computer vision can help to rapidly detect such outbreaks by automatically analyzing imagery from reefs for COTS abundance. In [50], a dataset of three videos from manta tows (a snorkeler with a camera, towed behind a boat) in Australia were annotated with bounding boxes of COTS. The videos contain 12347, 11374, and 10759 frames at 720×1280 px with 3065, 6384, and 2449 bounding-box annotations of COTS respectively. We use the bounding-boxes as segmentation labels to fine-tune a SegFormer-MiT-B5 pre-trained on Coralscapes for 10 epochs, with bounding boxes being drawn around connected components of the prediction for evaluation. The results (Table IV) show that pre-training on Coralscapes leads to a substantial improvement compared to using an ImageNet pre-trained backbone. Compared to a Yolo-V8-L baseline, the transfer-learned segmentation model consistently improves the mAP@50, but the mAP@50-95, which is more sensitive to the bounding box locations, is lower. This is likely due to the naive approach to transforming segmentations into bounding boxes.

Method	Label Style	Accuracy					mIoU				
		5	10	25	300	All	5	10	25	300	All
PLAS [64]	Grid	65.73	70.18	73.60	83.72	-	19.48	26.01	36.27	52.83	-
HIL-S (DinoV2+K-NN) [65]	Grid	72.24	77.64	85.93	85.93	-	25.66	34.80	43.41	54.07	-
HIL-S (DinoV2+K-NN) [65]	Human-in-the-Loop	74.53	71.04	81.69	86.29	-	32.96	38.21	45.46	54.62	-
DeeplabV3+ - ResNet50	Random	80.29	81.21	83.18	86.54	87.09	36.39	38.54	43.24	51.59	52.65
(Ours) DeeplabV3+ - ResNet50	Random	81.56	82.70	83.87	86.66	87.56	40.27	42.53	46.14	53.22	55.68
SegFormer - MiT-B2	Random	84.85	85.79	87.27	89.99	90.03	47.59	52.08	56.27	63.03	65.84
(Ours) SegFormer - MiT-B2	Random	84.98	86.18	87.44	89.04	89.92	48.36	52.41	55.80	62.01	65.42
SegFormer - MiT-B5 (Decoder)	Random	76.41	77.58	79.15	81.59	81.98	27.10	30.22	32.98	39.32	40.70
(Ours) SegFormer - MiT-B5 (Decoder)	Random	81.35	82.05	82.89	84.85	85.26	36.04	37.62	40.31	44.33	45.62

Table III: Transfer learning on UCSD Mosaics using 5 / 10 / 25 / 300 / all labeled pixels per image. Models marked with (Ours) are pre-trained on Coralscapes.

Method	Video 0		Video 1		Video 2	
	mAP	mAP	mAP	mAP	mAP	mAP
	@50	@50-95	@50	@50-95	@50	@50-95
YoloV8-L	0.496	0.295	0.384	0.226	0.413	0.222
SegFormer - MiT-B5	0.443	0.199	0.359	0.141	0.364	0.082
(Ours) SegFormer - MiT-B5	0.542	0.288	0.393	0.196	0.512	0.194

Table IV: COTS detection results in a leave-one-out test-set setting for each of the three videos.

Underwater 3D Mapping from Videos. Even though 3D photogrammetry is widely used to map coral reefs [16], [46], [74], [67], [15], commonly used pipelines are designed for in-air mapping and often struggle in underwater environments. The detrimental effects of the water column include changes in the color of objects relative to the distance to the camera and induced blur [4], [70]. Dynamic objects such as moving fish or divers also lead to artifacts in the 3D reconstructions or even complete failure. Therefore, image collections are usually highly curated [7] to exclude images with such objects or with the water column visible.

Instead of constraining the input data, general purpose semantic segmentation from Coralscapes can be used to mask out unwanted classes during 3D reconstruction pipelines, expanding the range of useful images for 3D reconstruction. We qualitatively demonstrate this on three reef videos, using GLOMAP [61] (at 2FPS in sequential matching mode) for sparse SfM, on which we run dense reconstruction with COLMAP [72] as well as 3D Gaussian Splatting (3DGS) [42] for novel view synthesis. Example images, masks and results are shown in Figure 5, where detrimental artifacts from the background are alleviated in dense reconstruction, and moving objects (fish and divers) are removed in 3DGS novel view synthesis. Similarly, masks obtained from models trained on Coralscapes can be used in other neural rendering techniques [47], [78], [48], or neural monocular SLAM systems [69].

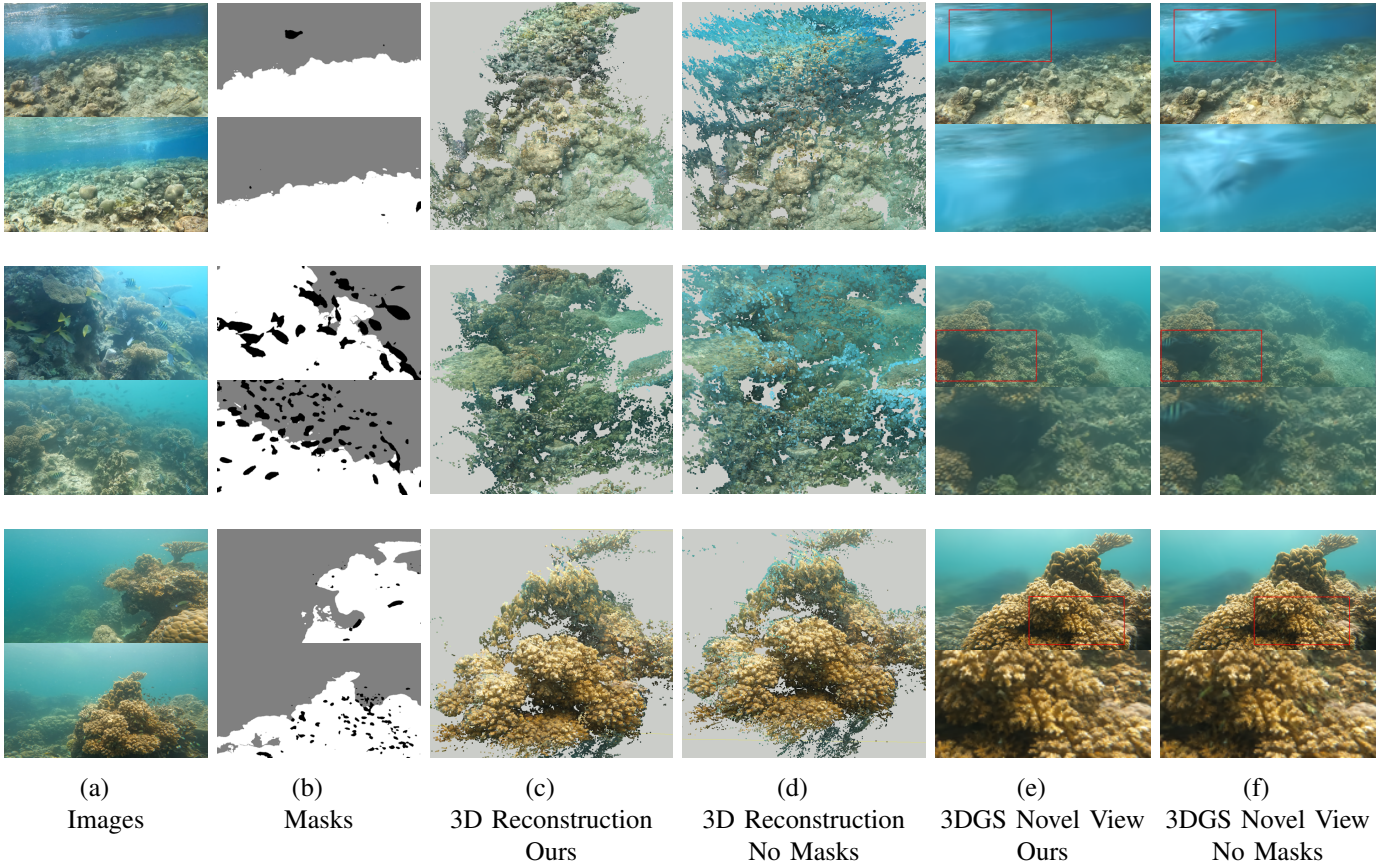


Fig. 5: Using masks obtained from Coralscapes to mask out unwanted classes can alleviate artifacts from dense 3D reconstruction and from novel view synthesis such as 3DGS.

VI. CONCLUSION AND FUTURE WORK

We present the Coralscapes dataset for semantic understanding in reef scenes, which is the first general purpose dataset of coral scenes that are densely and consistently annotated by trained experts. Coralscapes is representative of the challenges of semantic segmentation in reef scenes, including complex scenes with challenging environmental conditions. By benchmarking state-of-the-art computer vision models, we confirm that Coralscapes is a challenging benchmark due to the imbalance in class occurrence and feature details. The Coralscapes dataset aims to catalyze the development of computer vision applications in coral reef science and conservation. Future work will extend Coralscapes to further increase the diversity (e.g. new biogeographic regions, different camera models, onboard lighting), and add more fine-grained classes, also including fish genera or species, and include instance-level annotations for fish.

VII. ACKNOWLEDGMENTS

We thank Dr. Ibrahim Souleiman Abdallah (University of Djibouti, Djibouti), Prof. Osama S. Saad, Mustafa Altaib Mohammed, and Abualdrdaa Mirgani Hajbakhit Omer (Red Sea University, Port Sudan), as well as Dr. Zekeria Zekeria and Temesgen Gebremeskel (Mai-Nefhi College of Asmara, Eritrea) for helping to acquire video material. We thank Dr. Ali Al-Sawalmih and Tariq Al-Salman (Marine Science Station, Aqaba), Prof. Maoz Fine, Nahum Sela (InterUniversity Institute of Marine Science, Eilat), Dr. Assaf Zvuloni (Nature Reserve Authority Israel), the Aqaba Special Economic Zone Authority of Jordan, the Ministry of Environment and Sustainable Development of Djibouti and the University of Djibouti, the Ministry of Marine Resources of Eritrea and the Mai-Nefhi College of Asmara, Eritrea for their support for enabling us to collect the videos. We thank Chiara Freneix, Géza Soldati, Camille Perrin, Amélie Menoud, Ines Stiti, Martin Métier, and Antoine Carnal for help with the frame annotations. The main body of data in this study were collected in the framework of expeditions organized by Transnational Red Sea Center, which is hosted by the Laboratory for Biological Geochemistry at EPFL. This work was funded in part by FNS grant 205321_212614, as well as EPFL and the Transnational Red Sea Center.

REFERENCES

- [1] 2d semantic labeling contest: Potsdam. Online, Mar. 2020. Available: <https://www.isprs.org/education/benchmarks/UrbanSemLab/2d-sem-label-potsdam.aspx>. 2
- [2] 2d semantic labeling contest: Vaihingen. Online, Mar. 2020. Available: <https://www.isprs.org/education/benchmarks/UrbanSemLab/2d-sem-label-vaihingen.aspx>. 2
- [3] AIMS. Reefcloud.ai. <https://reefcloud.ai/>, 2023. 3
- [4] Derya Akkaynak and Tali Treibitz. Sea-thru: A method for removing water from underwater images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1682–1691, 2019. 9
- [5] Inigo Alonso, Matan Yuval, Gal Eyal, Tali Treibitz, and Ana C Murillo. Coralseg: Learning coral segmentation from sparse annotations. *Journal of Field Robotics*, 36(8):1456–1477, 2019. 2, 3, 4, 8
- [6] Mitchell B. Lyons, Chris M. Roelfsema, Emma V. Kennedy, Eva M. Kovacs, Rodney Borrego-Acevedo, Kathryn Markey, Meredith Roe, Dobby M. Yuwono, Daniel L. Harris, Stuart R. Phinn, et al. Mapping the world’s coral reefs using a global multiscale earth observation framework. *Remote Sensing in Ecology and Conservation*, 6(4):557–568, 2020. 2
- [7] Daniel TI Bayley and Andrew OM Mogg. A protocol for the large-scale analysis of reefs using structure from motion photogrammetry. *Methods in Ecology and Evolution*, 11(11):1410–1420, 2020. 9
- [8] Oscar Beijbom, Peter J. Edmunds, David I. Kline, B. Greg Mitchell, and David Kriegman. Automated annotation of coral reef survey images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Providence, Rhode Island, 2012. 2, 3
- [9] Oscar Beijbom, Peter J Edmunds, Chris Roelfsema, Jennifer Smith, David I Kline, Benjamin P Neal, Matthew J Dunlap, Vincent Moriarty, Tung-Yung Fan, Chih-Jui Tan, et al. Towards automated annotation of benthic survey images: Variability of human experts and operational modes of automation. *PloS one*, 10(7):e0130312, 2015. 2, 3
- [10] Oscar Beijbom, Tali Treibitz, David I Kline, Gal Eyal, Adi Khen, Benjamin Neal, Yossi Loya, B Greg Mitchell, and David Kriegman. Improving automated annotation of benthic survey images using wide-band fluorescence. *Scientific reports*, 6(1):23166, 2016. 3
- [11] Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE transactions on medical imaging*, 37(11):2514–2525, 2018. 2
- [12] Michael Bewley, Ariell Friedman, Renata Ferrari, Nicole Hill, Renae Hovey, Neville Barrett, Ezequiel M Marzinelli, Oscar Pizarro, Will Figueira, Lisa Meyer, et al. Australian sea-floor survey data, with images and expert annotations. *Scientific data*, 2(1):1–13, 2015. 3
- [13] Hawthorne L Beyer, Emma V Kennedy, Maria Beger, Chaolun Allen Chen, Joshua E Cinner, Emily S Darling, C Mark Eakin, Ruth D Gates, Scott F Heron, Nancy Knowlton, et al. Risk-sensitive planning for conserving coral reefs under rapid climate change. *Conservation Letters*, 11(6):e12587, 2018. 1

- [14] Patrick Bilic, Patrick Christ, Hongwei Bran Li, Eugene Vorontsov, Avi Ben-Cohen, Georgios Kaissis, Adi Szeskin, Colin Jacobs, Gabriel Efrain Humpire Mamani, Gabriel Chartrand, et al. The liver tumor segmentation benchmark (lits). *Medical image analysis*, 84:102680, 2023. [2](#)
- [15] Pim Bongaerts, Caroline E Dubé, Katharine E Prata, Johanna C Gijsbers, Michelle Achlatis, and Alejandra Hernandez-Agreda. Reefscape genomics: leveraging advances in 3d imaging to assess fine-scale patterns of genomic variation on coral reefs. *Frontiers in Marine Science*, page 875, 2021. [9](#)
- [16] JHR Burns, D Delparte, RD Gates, and M Takabayashi. Integrating structure-from-motion photogrammetry with geospatial software as a novel technique for quantifying 3d ecological characteristics of coral reefs. *PeerJ*, 3:e1077, 2015. [9](#)
- [17] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018. [2](#)
- [18] Elisa Casella, Antoine Collin, Daniel Harris, Sebastian Ferse, Sonia Bejarano, Valeriano Parravicini, James L Hench, and Alessio Rovere. Mapping coral reefs using consumer-grade drones and structure from motion photogrammetry techniques. *Coral Reefs*, 36:269–275, 2017. [2](#)
- [19] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. [2](#)
- [20] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. [2](#)
- [21] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. [2](#), [6](#), [24](#)
- [22] Qimin Chen, Oscar Beijbom, Stephen Chan, Jessica Bouwmeester, and David Kriegman. A new deep learning engine for coralnet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3693–3702, 2021. [3](#)
- [23] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. [2](#)
- [24] CVAT.ai. Cvat doi 10.5281/zenodo.4009388. <https://github.com/cvat-ai/cvat>, 2023. [5](#)
- [25] Ilke Demir, Krzysztof Koperski, David Lindenbaum, Guan Pang, Jing Huang, Saikat Basu, Forest Hughes, Devis Tuia, and Ramesh Raskar. Deepglobe 2018: A challenge to parse the earth through satellite images. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 172–181, 2018. [2](#)
- [26] Adele M Dixon, Piers M Forster, Scott F Heron, Anne MK Stoner, and Maria Beger. Future loss of local-scale thermal refugia in coral reef ecosystems. *Plos Climate*, 1(2):e0000004, 2022. [1](#)
- [27] Clinton B Edwards, Yoan Eynaud, Gareth J Williams, Nicole E Pedersen, Brian J Zgliczynski, Arthur CR Gleason, Jennifer E Smith, and Stuart A Sandin. Large-area imaging reveals biologically driven non-random spatial patterns of corals at a remote reef. *Coral Reefs*, 36:1291–1305, 2017. [2](#), [3](#), [8](#)
- [28] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111:98–136, 2015. [2](#)
- [29] Maoz Fine, Hezi Gildor, and Amatzia Genin. A coral reef refuge in the red sea. *Global change biology*, 19(12):3640–3647, 2013. [1](#)
- [30] Rebecca Fisher, Rebecca A O’Leary, Samantha Low-Choy, Kerrie Mengersen, Nancy Knowlton, Russell E Brainard, and M Julian Caley. Species richness on coral reefs and the pursuit of convergent global estimates. *Current Biology*, 25(4):500–505, 2015. [1](#)
- [31] Daniel P Furtado, Edson A Vieira, Wildna Fernandes Nascimento, Kelly Y Inagaki, Jessica Bleuel, Marco Antonio Zanata Alves, Guilherme O Longo, and Luiz S Oliveira. # deolhonoscorais: a polygonal annotated dataset to optimize coral monitoring. *PeerJ*, 11:e16219, 2023. [3](#)
- [32] Manuel González-Rivero, Alberto Rodríguez-Ramírez, Oscar Beijbom, Peter Dalton, Emma V Kennedy, Benjamin P Neal, Julie Vercelloni, Pim Bongaerts, Anjani Ganase, Dominic EP Bryant, et al. Seaview survey photo-quadrat and image classification dataset. 2019. [3](#)
- [33] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364, 2019. [2](#)
- [34] Karlo Hock, Nicholas H Wolff, Scott A Condie, Kenneth RN Anthony, and Peter J Mumby. Connectivity networks reveal the risks of crown-of-thorns starfish outbreaks on the great barrier reef. *Journal of applied ecology*, 51(5):1188–1196, 2014. [8](#)
- [35] Brian M Hopkinson, Andrew C King, Daniel P Owen, Matthew Johnson-Roberson, Matthew H Long, and Suchendra M Bhandarkar. Automated classification of three-dimensional reconstructions of coral reefs using convolutional neural networks. *PloS one*, 15(3):e0230671, 2020. [3](#)
- [36] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. [6](#)
- [37] Pavel Iakubovskii. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2019. [24](#)
- [38] Md Jahidul Islam, Chelsey Edge, Yuyang Xiao, Peigen Luo, Muntaqim Mehtaz, Christopher Morse, Sadman Sakib Enan, and Junaed Sattar. Semantic segmentation of underwater imagery: Dataset and benchmark. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1769–1776. IEEE, 2020. [3](#)
- [39] Sonain Jamil, MuhibUr Rahman, and Amir Haider. Bag of features (bof) based deep learning framework for bleached corals detection. *Big Data and Cognitive Computing*, 5(4):53, 2021. [2](#), [3](#)
- [40] Glenn Jocher, Jing Qiu, and Ayush Chaurasia. Ultralytics, 01 2023. [26](#)
- [41] Keivan Kabiri, Hamid Rezai, and Masoud Moradi. A drone-based method for mapping the coral reefs in the shallow coastal waters—case study: Kish island, persian gulf. *Earth Science Informatics*, 13(4):1265–1274, 2020. [2](#)
- [42] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. [9](#)

- [43] Andrew King, Suchendra M Bhandarkar, and Brian M Hopkinson. A comparison of deep learning methods for semantic segmentation of coral reef survey images. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 1394–1402, 2018. 3
- [44] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 3, 5
- [45] Thomas Krueger, Noa Horwitz, Julia Bodin, Maria-Evangelia Giovani, Stéphane Escrig, Anders Meibom, and Maoz Fine. Common reef-building coral in the northern red sea resistant to elevated temperature and acidification. *Royal Society open science*, 4(5):170038, 2017. 1
- [46] Javier X Leon, Chris M Roelfsema, Megan I Saunders, and Stuart R Phinn. Measuring coral reef terrain roughness using ‘structure-from-motion’ close-range photogrammetry. *Geomorphology*, 242:21–28, 2015. 9
- [47] Deborah Levy, Amit Peleg, Naama Pearl, Dan Rosenbaum, Derya Akkaynak, Simon Korman, and Tali Treibitz. Seathru-nerf: Neural radiance fields in scattering media. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 56–65, 2023. 9
- [48] Huapeng Li, Wenxuan Song, Tianao Xu, Alexandre Elsig, and Jonas Kulhanek. Watersplating: Fast underwater 3d scene reconstruction using gaussian splatting. *arXiv preprint arXiv:2408.08206*, 2024. 9
- [49] Jiwei Li, David E Knapp, Nicholas S Fabina, Emma V Kennedy, Kirk Larsen, Mitchell B Lyons, Nicholas J Murray, Stuart R Phinn, Chris M Roelfsema, and Gregory P Asner. A global coral reef probability map generated using convolutional neural networks. *Coral Reefs*, 39:1805–1815, 2020. 2
- [50] Jiajun Liu, Brano Kusy, Ross Marchant, Brendan Do, Torsten Merz, Joey Crosswell, Andy Steven, Nic Heaney, Karl von Richter, Lachlan Tychsen-Smith, et al. The csiro crown-of-thorn starfish detection dataset. *arXiv preprint arXiv:2111.14311*, 2021. 8, 26
- [51] Scott C Lowe, Benjamin Misiuk, Isaac Xu, Shakhboz Abdulazizov, Amit R Baroi, Alex C Bastos, Merlin Best, Vicki Ferrini, Ariell Friedman, Deborah Hart, et al. Benthicnet: A global compilation of seafloor images for deep learning applications. *arXiv preprint arXiv:2405.05241*, 2024. 3
- [52] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022. 24
- [53] Valérie Masson-Delmotte, Panmao Zhai, Anna Pirani, Sarah L Connors, Clotilde Péan, S Berger, N Caud, Y Chen, L Goldfarb, MI Gomis, et al. Climate change 2021: the physical science basis. *Contribution of working group I to the sixth assessment report of the intergovernmental panel on climate change*, page 2, 2021. 1
- [54] Valérie Masson-Delmotte, Panmao Zhai, Hans-Otto Pörtner, Debra Roberts, Jim Skea, Priyadarshi R Shukla, Anna Pirani, Wilfran Moufouma-Okia, Clotilde Péan, Roz Pidcock, et al. Global warming of 1.5 c. *An IPCC Special Report on the impacts of global warming of*, 1(5), 2018. 1
- [55] Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging*, 34(10):1993–2024, 2014. 2
- [56] MERMAID. Marine ecological research management aid (mermaid). <https://datamermaid.org/>, 2023. 3
- [57] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Buló, and Peter Kotschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *Proceedings of the IEEE international conference on computer vision*, pages 4990–4999, 2017. 2
- [58] NOAA. Fact sheet: Coral reefs. <https://www.coast.noaa.gov/states/fast-facts/coral-reefs.html>, 2022. Accessed: 2022-09-09. 1
- [59] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 6, 24
- [60] Eslam O Osman, David J Smith, Maren Ziegler, Benjamin Kürten, Constanze Conrad, Khaled M El-Haddad, Christian R Voolstra, and David J Suggett. Thermal refugia against coral bleaching throughout the northern red sea. *Global change biology*, 24(2):e474–e484, 2018. 1
- [61] Linfei Pan, Dániel Baráth, Marc Pollefeys, and Johannes L Schönberger. Global structure-from-motion revisited. In *European Conference on Computer Vision*, pages 58–77. Springer, 2025. 9
- [62] Morgan S Pratchett, Ciemon F Caballes, Jairo A Rivera-Posada, and Hugh PA Sweatman. Limits to understanding and managing outbreaks of crown-of-thorns starfish (*acanthaster* spp.). In *Oceanography and Marine Biology*, pages 133–200. CRC Press, 2014. 8
- [63] Morgan S Pratchett, Ciemon F Caballes, Jennifer C Wilmes, Samuel Matthews, Camille Mellin, Hugh PA Sweatman, Lauren E Nadler, Jon Brodie, Cassandra A Thompson, Jessica Hoey, et al. Thirty years of research on crown-of-thorns starfish (1986–2016): scientific advances and emerging opportunities. *Diversity*, 9(4):41, 2017. 8
- [64] Scarlett Raine, Ross Marchant, Brano Kusy, Frederic Maire, and Tobias Fischer. Point label aware superpixels for multi-species segmentation of underwater imagery. *IEEE Robotics and Automation Letters*, 7(3):8291–8298, 2022. 3, 8
- [65] Scarlett Raine, Ross Marchant, Brano Kusy, Frederic Maire, Niko Sunderhauf, and Tobias Fischer. Human-in-the-loop segmentation of multi-species coral imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2723–2732, 2024. 8
- [66] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–12188, 2021. 6, 24
- [67] Vincent Raoult, Peter A David, Sally F Dupont, Ciaran P Mathewson, Samuel J O’Neill, Nicholas N Powell, and Jane E Williamson. GoproTM as an underwater photogrammetry tool for citizen science. *PeerJ*, 4:e1960, 2016. 9
- [68] Alina Raphael, Zvy Dubinsky, David Iluz, Jennifer IC Benichou, and Nathan S Netanyahu. Deep neural network recognition of shallow water corals in the gulf of eilat (aqaba). *Scientific reports*, 10(1):12959, 2020. 2, 3
- [69] Jonathan Sauder, Guilhem Banc-Prandi, Anders Meibom, and Devis Tuia. Scalable semantic 3d mapping of coral reefs with deep learning. *Methods in Ecology and Evolution*, 15(5):916–934, 2024. 9

- [70] Jonathan Sauder and Devis Tuia. Self-supervised underwater caustics removal and descattering via deep monocular slam. In *European Conference on Computer Vision*, pages 214–232. Springer, 2024. [9](#)
- [71] Romain Savary, Daniel J Barshis, Christian R Voolstra, Anny Cárdenas, Nicolas R Evensen, Guilhem Banc-Prandi, Maoz Fine, and Anders Meibom. Fast and pervasive transcriptomic resilience and acclimation of extremely heat-tolerant coral holobionts from the northern red sea. *Proceedings of the National Academy of Sciences*, 118(19):e2023298118, 2021. [1](#)
- [72] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. [9](#)
- [73] ASM Shihavuddin. Coral reef dataset, v2. [3](#)
- [74] Curt D Storlazzi, Peter Dartnell, Gerald A Hatcher, and Ann E Gibbs. End of the chain? rugosity and fine-scale bathymetry from existing underwater digital imagery using structure-from-motion (sfm) technology. *Coral Reefs*, 35(3):889–894, 2016. [9](#)
- [75] Juan Terven, Diana-Margarita Córdova-Esparza, and Julio-Alejandro Romero-González. A comprehensive review of yolo architectures in computer vision: From yolov1 to yolov8 and yolo-nas. *Machine learning and knowledge extraction*, 5(4):1680–1716, 2023. [26](#)
- [76] Peter Todd. Todd p. a. — morphological plasticity in scleractinian corals. biological reviews. *Biological reviews of the Cambridge Philosophical Society*, 83:315–37, 09 2008. [4](#)
- [77] Christian R Voolstra, Jacob J Valenzuela, Serdar Turkarslan, Anny Cárdenas, Benjamin CC Hume, Gabriela Perna, Carol Buitrago-López, Katherine Rowe, Monica V Orellana, Nitin S Baliga, et al. Contrasting heat stress response patterns of coral holobionts across the red sea suggest distinct mechanisms of thermal tolerance. *Molecular ecology*, 30(18):4466–4480, 2021. [1](#)
- [78] Haoran Wang, Nantheera Anantrasirichai, Fan Zhang, and David Bull. Uw-gs: Distractor-aware 3d gaussian splatting for enhanced underwater scene reconstruction. *arXiv preprint arXiv:2410.01517*, 2024. [9](#)
- [79] David A Westcott, Cameron S Fletcher, Frederieke J Kroon, Russell C Babcock, Eva E Plagányi, Morgan S Pratchett, and Mary C Bonin. Relative efficacy of three approaches to mitigate crown-of-thorns starfish outbreaks on australia’s great barrier reef. *Scientific reports*, 10(1):12594, 2020. [8](#)
- [80] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019. [24](#)
- [81] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34:12077–12090, 2021. [6](#), [24](#)
- [82] Matan Yuval, Iñigo Alonso, Gal Eyal, Dan Tchernov, Yossi Loya, Ana C Murillo, and Tali Treibitz. Repeatable semantic reef-mapping through photogrammetry and label-augmentation. *Remote Sensing*, 13(4):659, 2021. [3](#)
- [83] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. [2](#)
- [84] Ziqiang Zheng, Haixin Liang, Binh-Son Hua, Yue Him Wong, Put Ang, Apple Pui Yi Chui, and Sai-Kit Yeung. Coralscop: Segment any coral image on this planet. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28170–28180, 2024. [3](#), [5](#)
- [85] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302–321, 2019. [2](#)
- [86] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep learning in medical image analysis and multimodal learning for clinical decision support: 4th international workshop, DLMIA 2018, and 8th international workshop, ML-CDS 2018, held in conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, proceedings 4*, pages 3–11. Springer, 2018. [6](#), [24](#)

APPENDIX

A. Appendix: Additional Dataset Information

Here, we give further information about the Coralscapes dataset and the annotated classes. We show the number of annotated pixels per class in Fig. 6, and the image frequency of each classes (on what percentage of images this class is present) in Fig. 7.

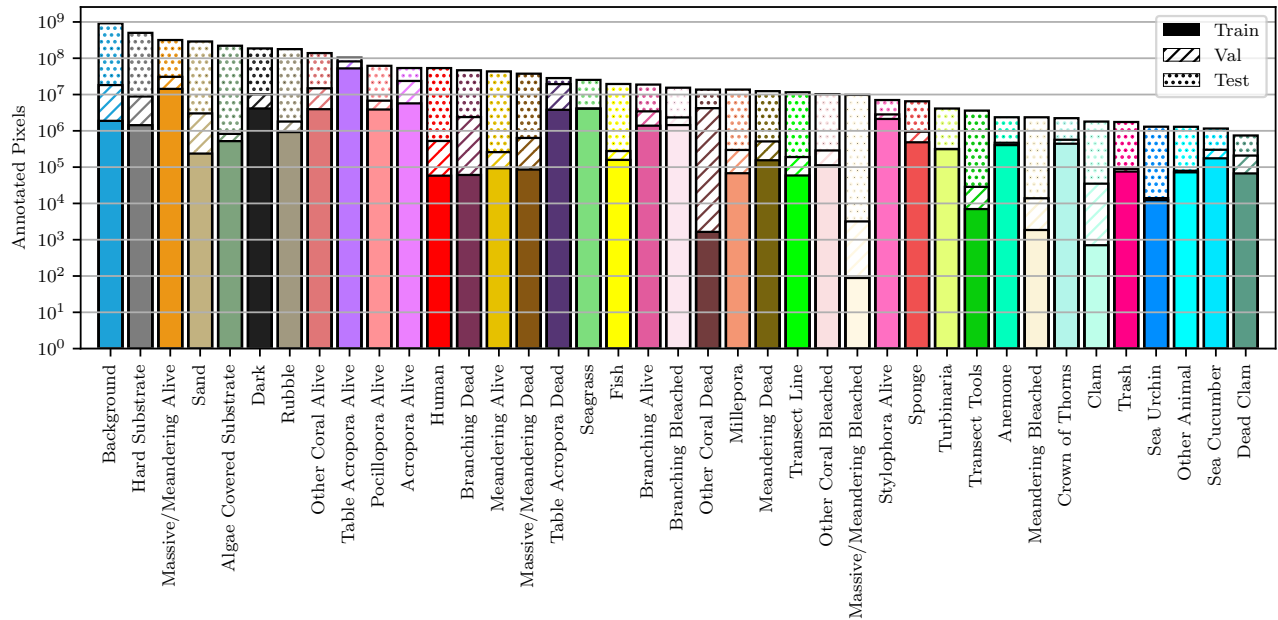


Fig. 6: Number of annotated pixels by class on a log-scale (with linear proportions for train/val/test).

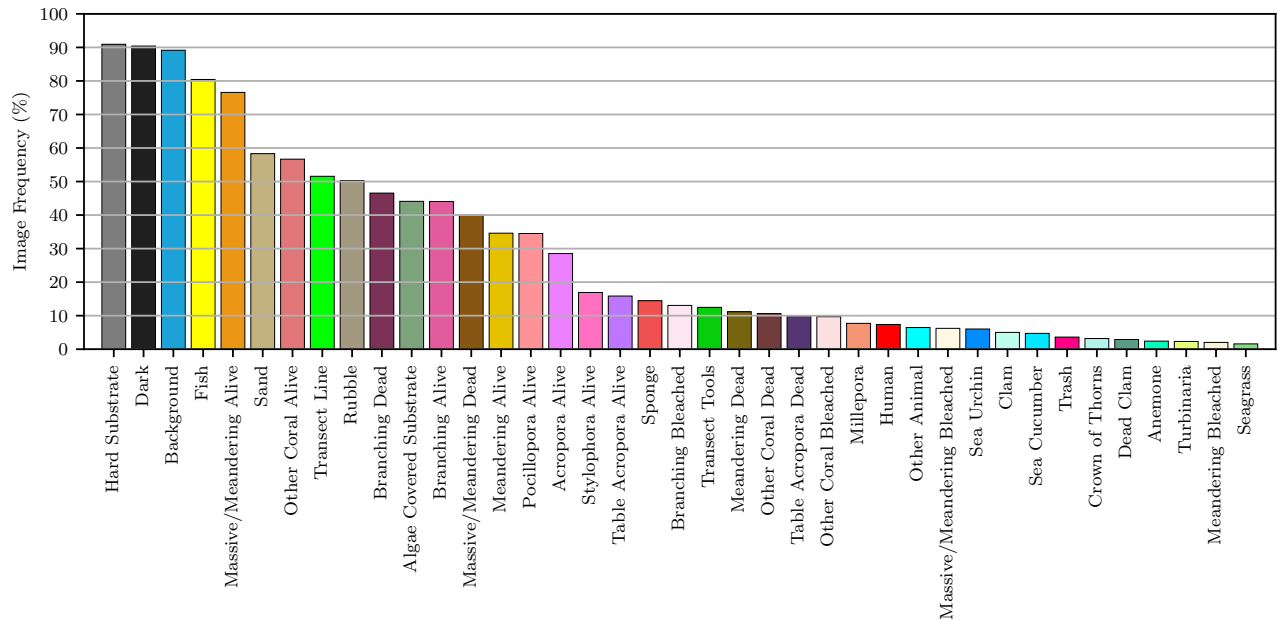


Fig. 7: Percentage of images with each class present.

Choosing the best resolution for a semantic segmentation dataset is challenging, as the label set should capture all classes of interest while at the same time having enough samples of each class to be able to accurately assess the performance of a machine learning system. In Figure 8, we propose a hierarchical scheme for summarizing the label set into coarser classes.

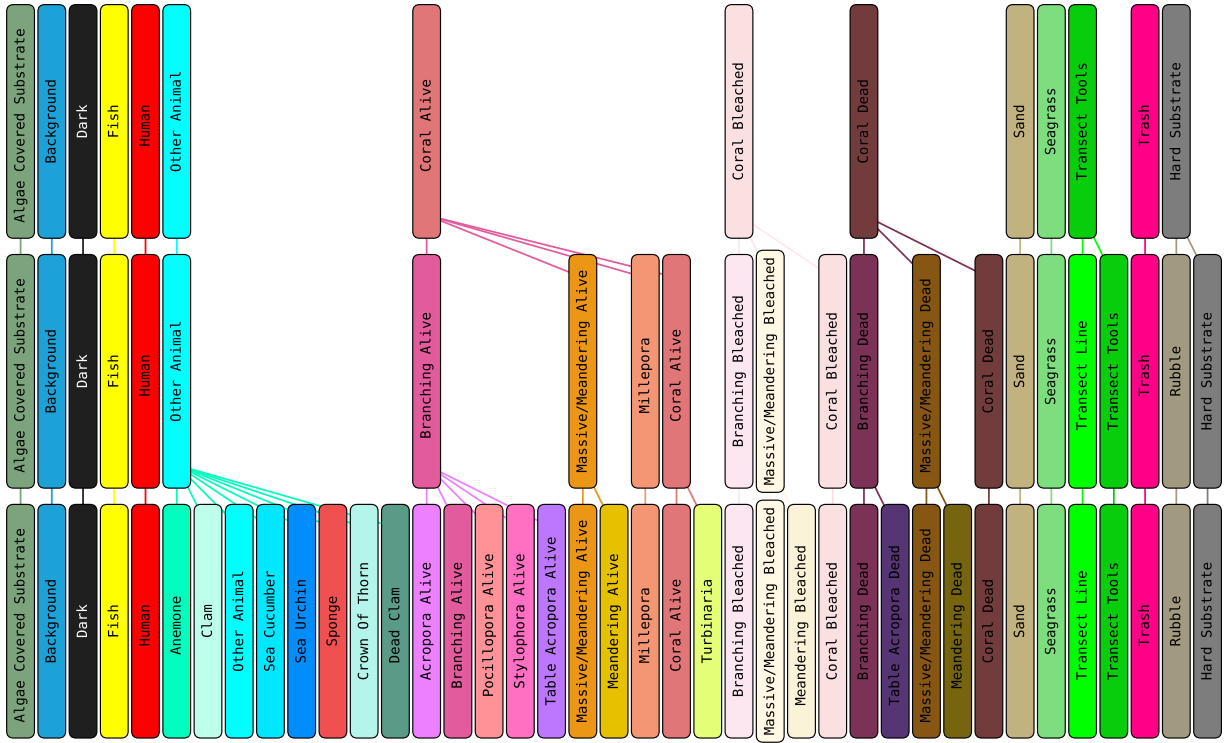

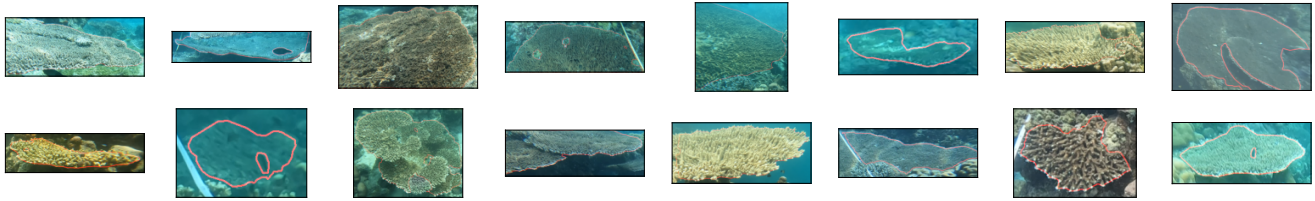



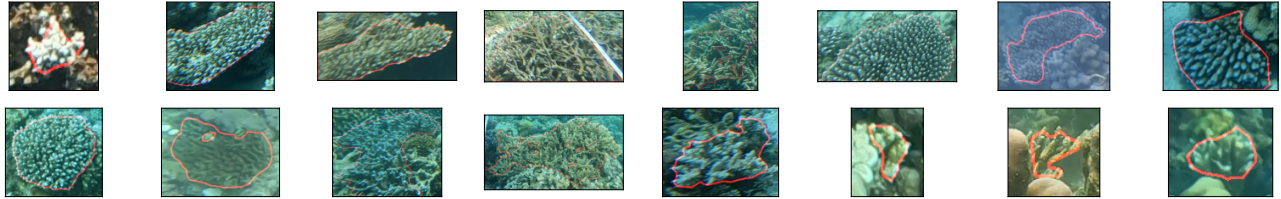
Fig. 8: Proposed label hierarchy for summarization.

B. Annotation Class Guide

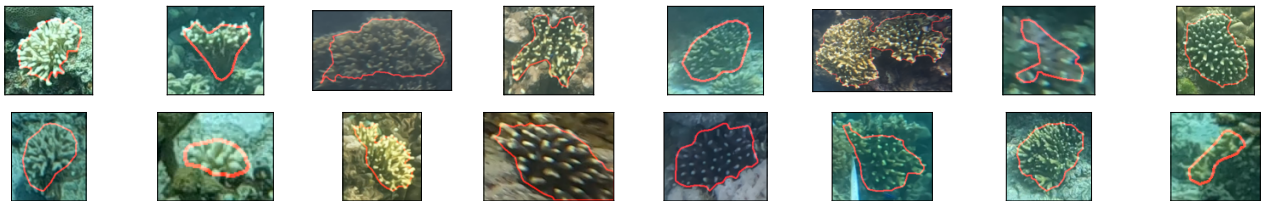
Table Acropora Alive  Acropora that grow in a tabular growth form.




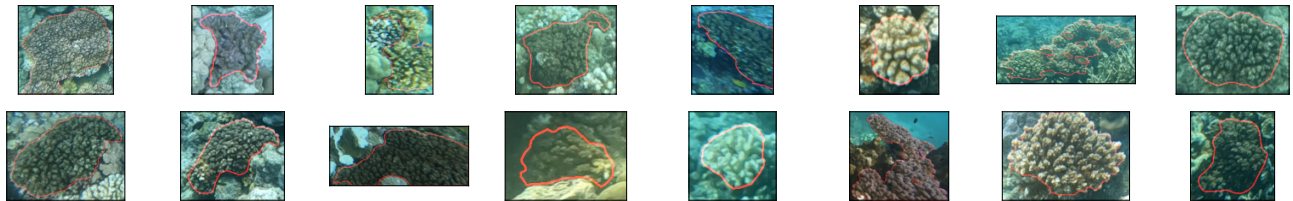
Acropora Alive  Acropora that do not grow in tabular growth form.




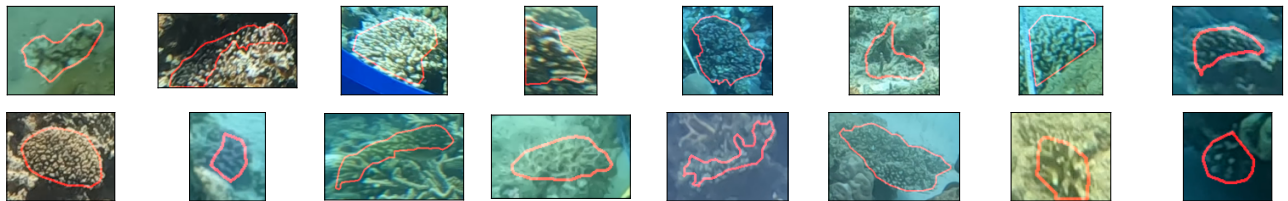
Stylophora Alive  Clearly identifiable Stylophora.



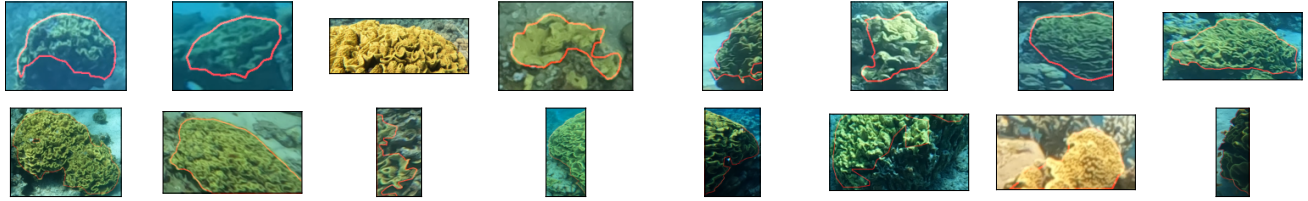
Pocillopora Alive  Clearly identifiable Pocillopora.




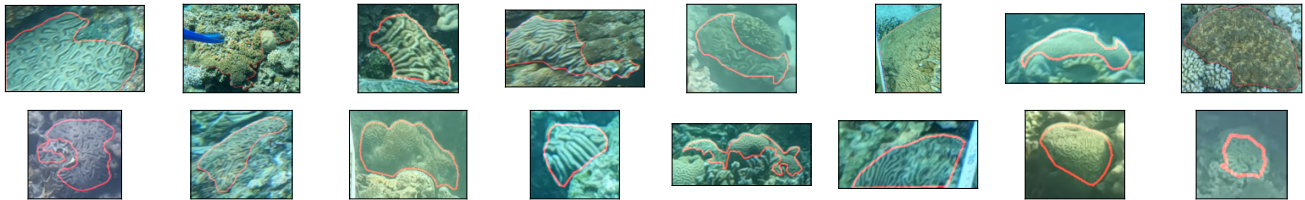
Branching Alive  Branching corals that can surely be determined to be alive, but do not fit in to 'table acropora alive', 'acropora', 'stylophora', or 'pocillopora' because they are from a different genus or appear slightly blurred in the context.




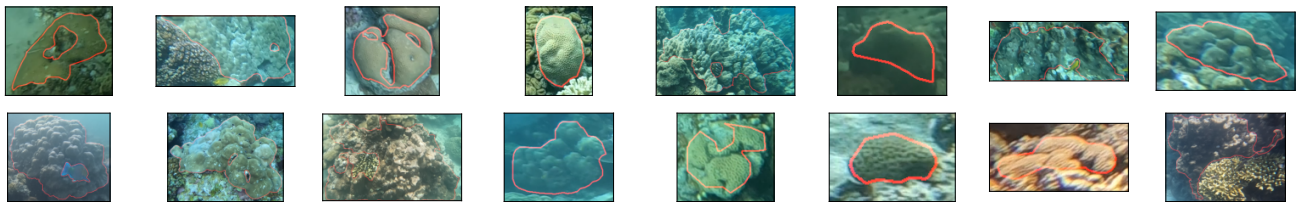
Turbinaria  Colloqually the ‘scroll’ coral. Does not include the macroalgae genus Turbinaria.




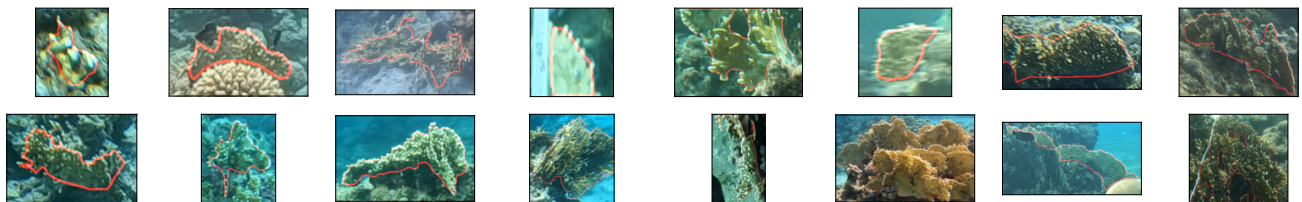
Meandering Alive  Corals of a meandering growth form. Includes *Platygyra*, *Lobophyllia*, *Symphyllia*.




Massive/Meandering Alive  Corals in a massive growth form. Prominently includes *Porites*, *Favia*, *Favites*, and many others. Includes some likely meandering corals (like *Platygyra*) where the meandering structure can not be clearly identified.



Millepora  The ‘fire coral’ *Millepora* is technically not a coral, but a hydrozoan. Appears most commonly in a branching form (*Millepora Dichotoma*).



Other Coral Alive  All other live corals. Includes corals in thin plate or encrusting growth form, soft corals, and corals that can not be clearly classified into the other classes.

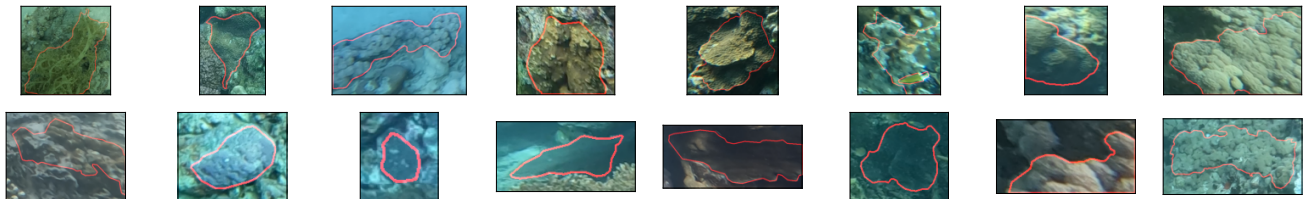

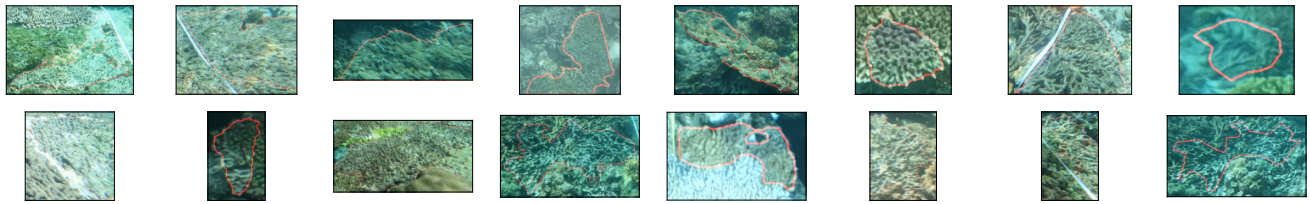

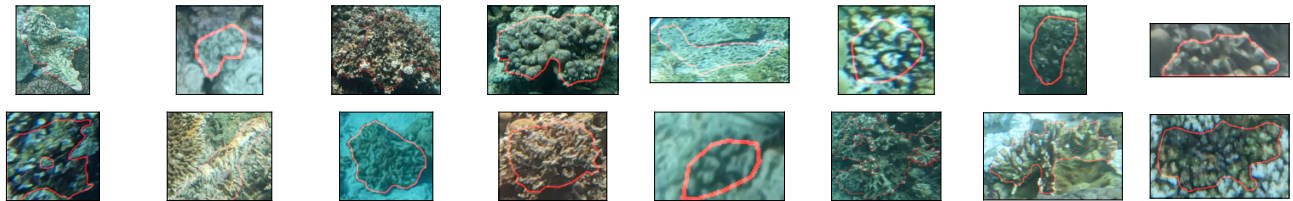


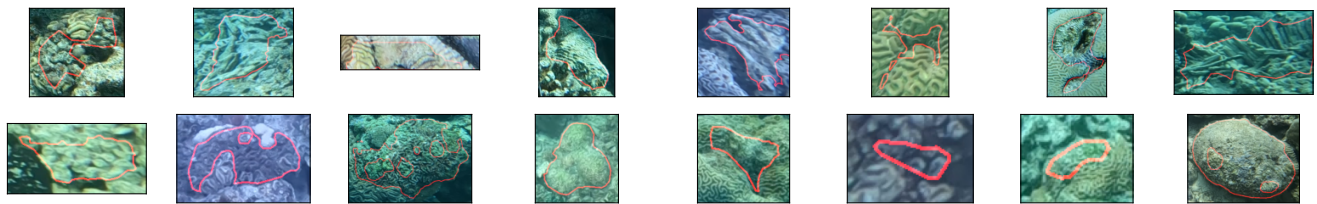
Table Acropora Dead  Dead Acropora tables which are dead or collapsed but are not yet overgrown by algae. Parts visibly overgrown by algae are labeled as ‘algae covered substrate’.




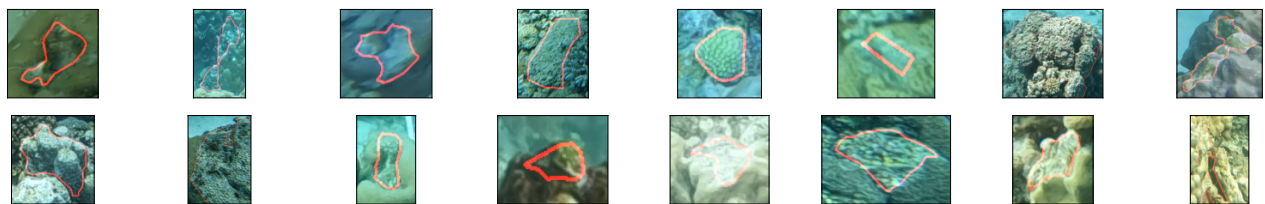
Branching Dead  Other dead branching coral, including acropora, pocillopora, and stylophora. When overgrown by algae, labeled as ‘algae covered substrate’.



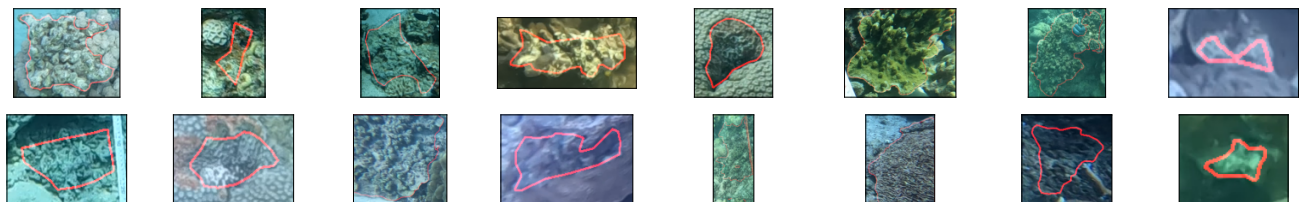
Meandering Dead  Dead meandering corals (Lobophyllia, Symphyllia, Platygyra, etc.).



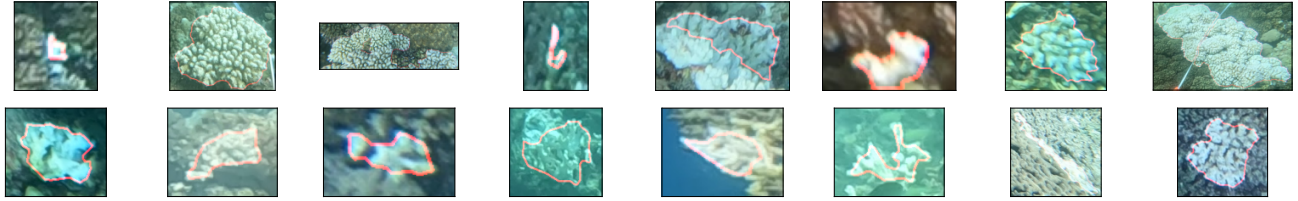
Massive/Meandering Dead  Dead massive corals or meandering corals that have decomposed enough so that the meandering structure is no longer visible.




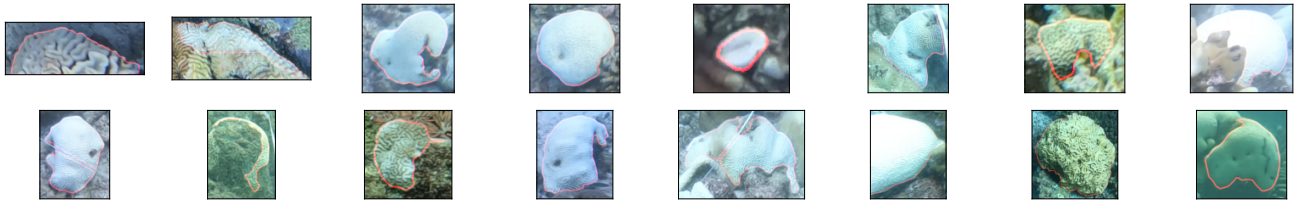
Other Coral Dead  All other dead coral skeletons.

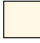


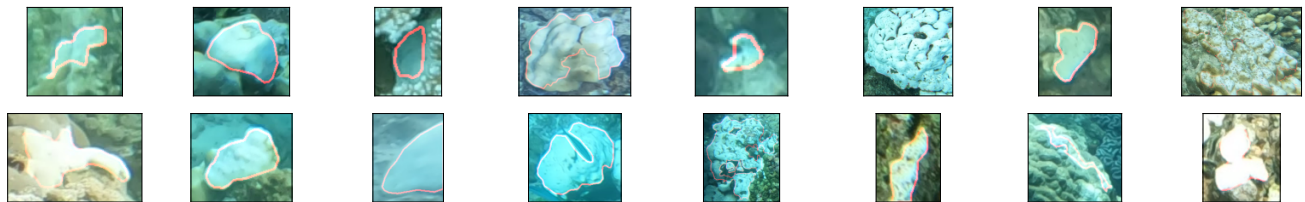
Branching Bleached  Bleached branching coral of all kinds, including bleached table acropora.



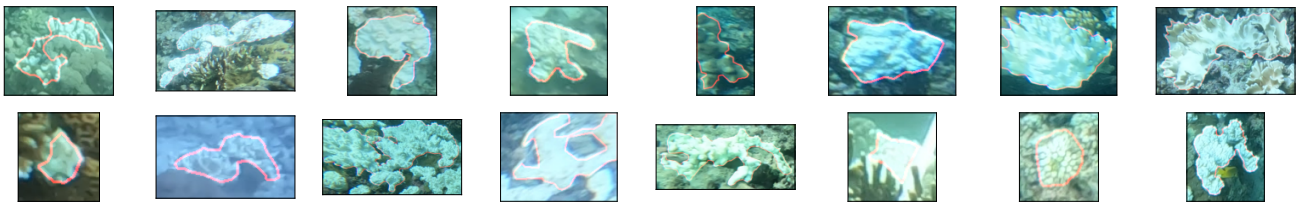
Meandering Bleached  Bleached meandering corals (Lobophyllia, Symphyllia, Platygyra, etc.).



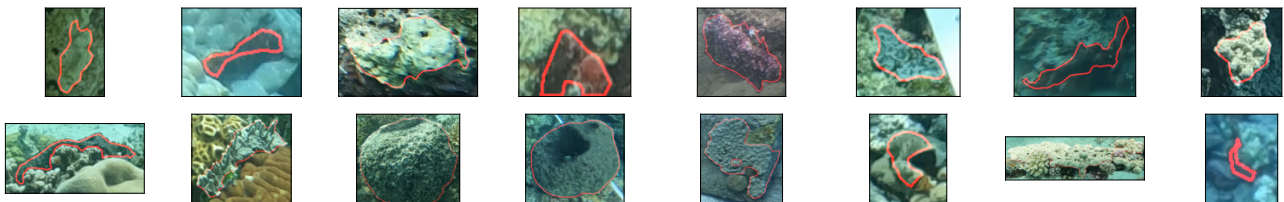
Massive/Meandering Bleached  Bleached massive corals or meandering corals that have decomposed enough so that the meandering structure is no longer visible.




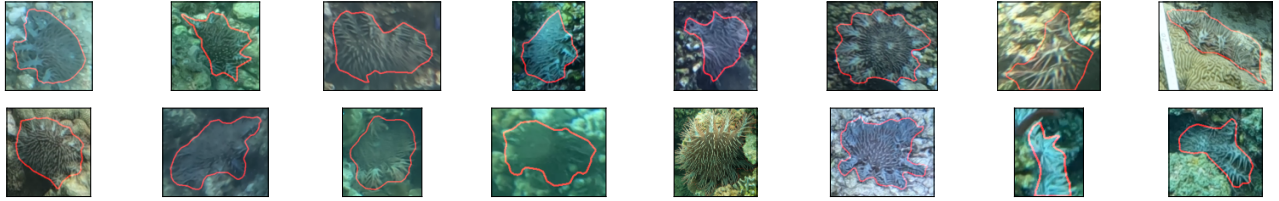
Other Coral Bleached  All other bleached coral, including bleached soft coral.



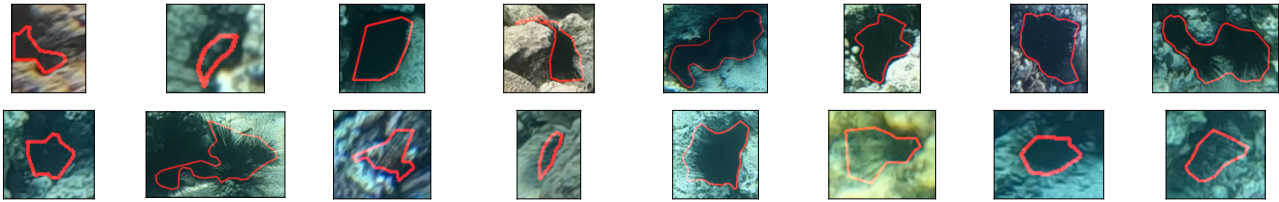
Sponge  Sponges of all kind.



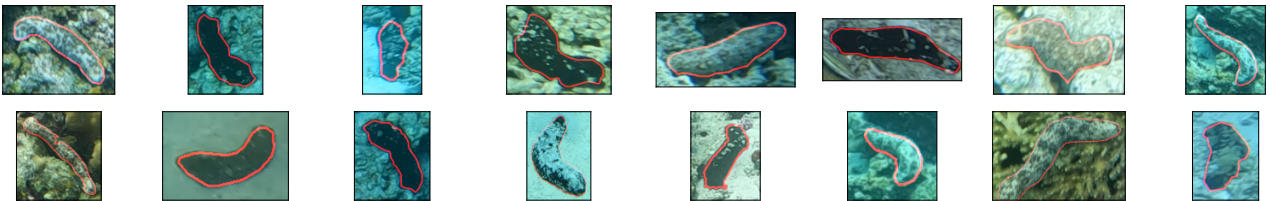
Crown of Thorns  *Acanthaster planci*, known to cause outbreaks in which it can severely reduce coral cover.



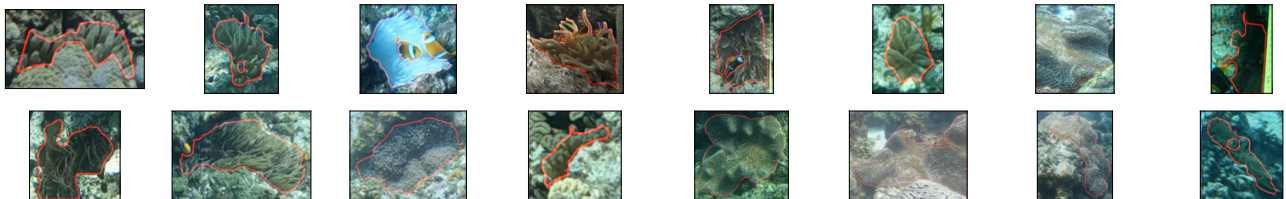
Sea Urchin 




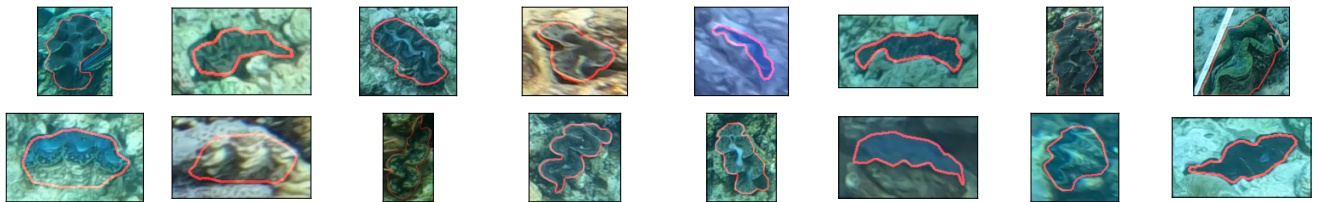
Sea Cucumber 



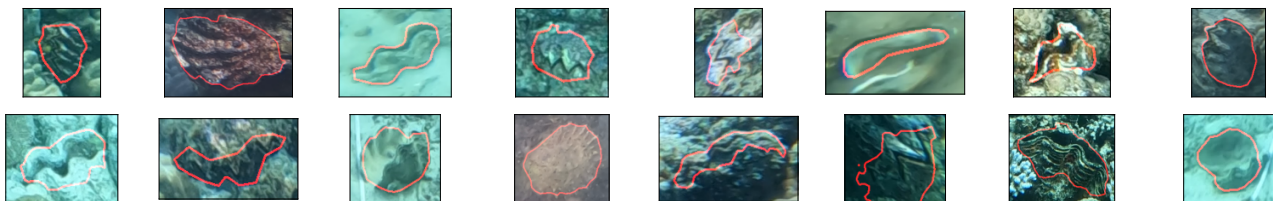
Anemone 

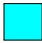


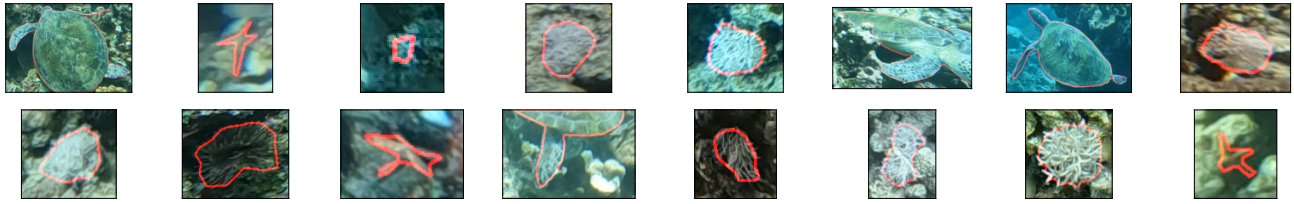
Clam  Live giant clams.




Dead Clam  Dead giant clams.




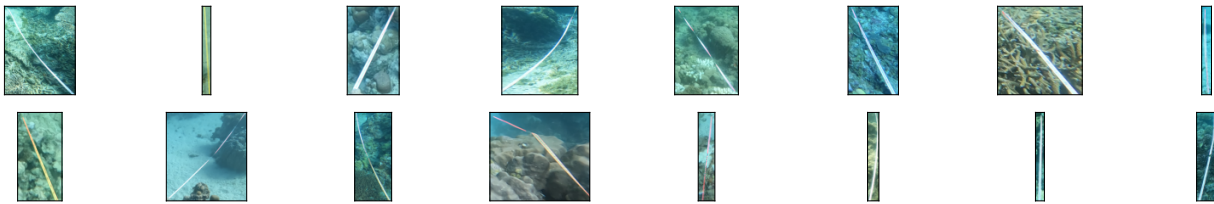
Other Animal  Includes starfish (except the crown-of-thorns starfish), feather worms, sea turtles, and other non-identifiable invertebrates or animals of which there are not enough annotations to warrant a separate class.




Trash  Includes all kinds of marine litter. Most common are plastic items including bags, cups, and bottles, aluminum cans and glass bottles, as well as abandoned fishing material, and parts of boats and machines.




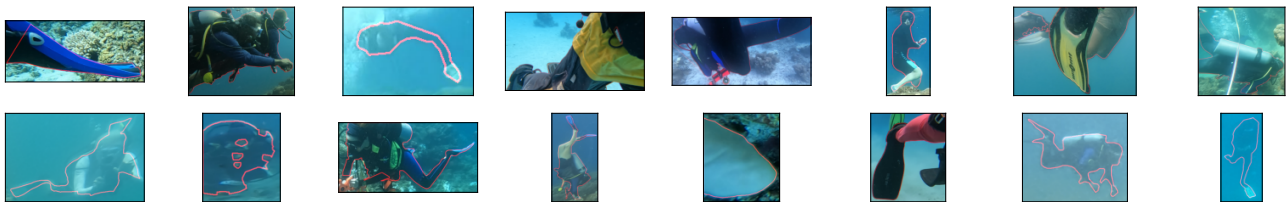
Transect Line  Rolled out transect tape. Excludes reels or other strands of rope that are laid out.




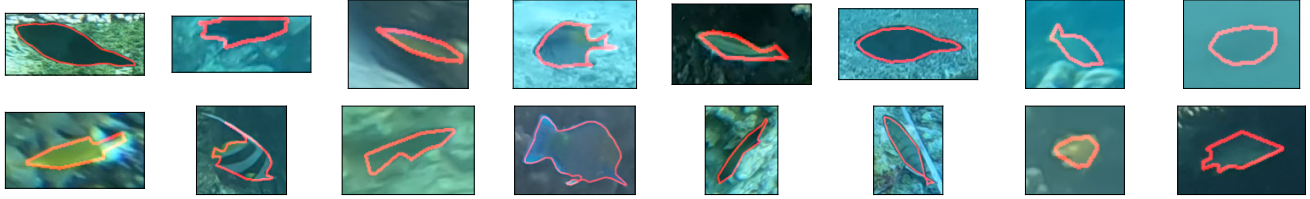
Transect Tools  This includes transect reels & spools, tags and markers placed on the reef, diving weights, surface marker buoys and the string attaching them to the ground or weights on the ground.




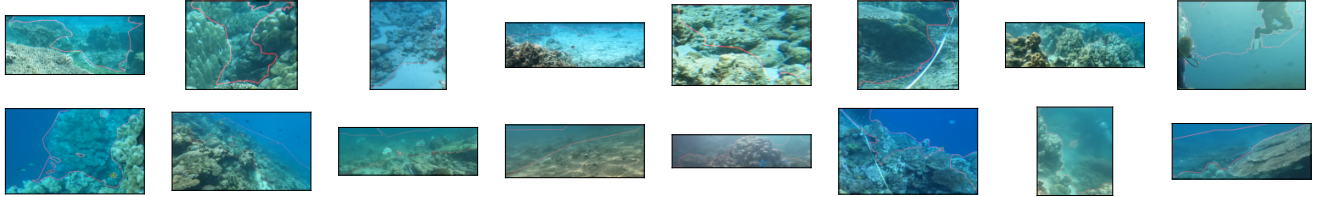
Human  The human class includes divers and snorkelers, of which sometimes only a hand or a fin is in the frame. Some ambiguity can arise when a human carries a transect tool that is not laid out on the ground, like a transect reel. In these cases, we generally decide this tool becomes part of the human polygon.




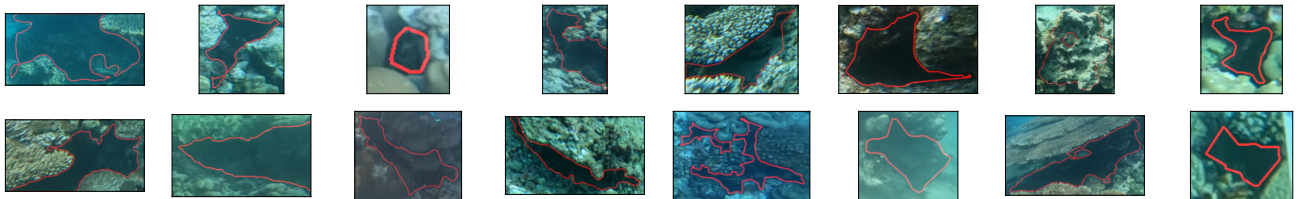
Fish  Fish of all kinds.




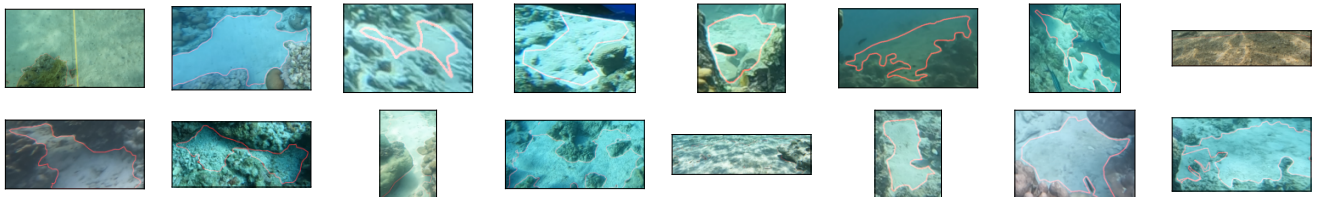
Background  This class is assigned to pixels that are too far away or too blurry to be classifiable into any other class. This includes the water surface, which is sometimes visible.




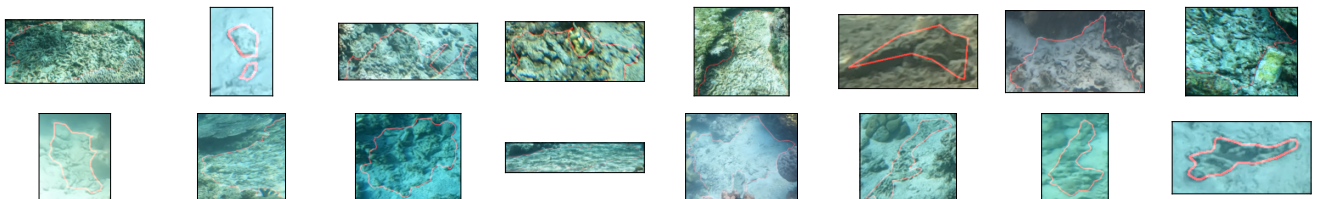
Dark  Parts of the image that are too dark to discern the benthic class.




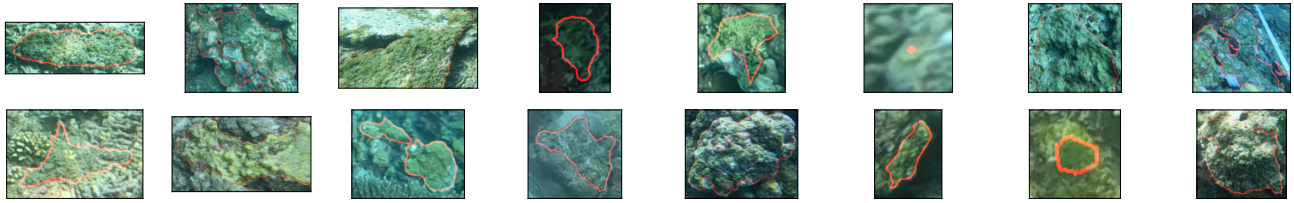
Sand  Loose, fine sand.




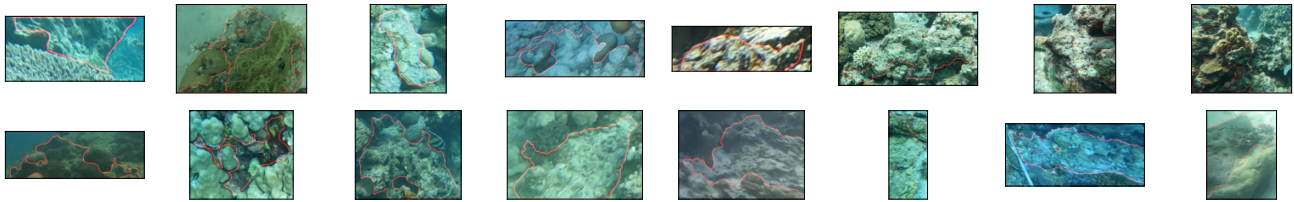
Rubble  Small loose fragments of rocky substrate or dead coral.



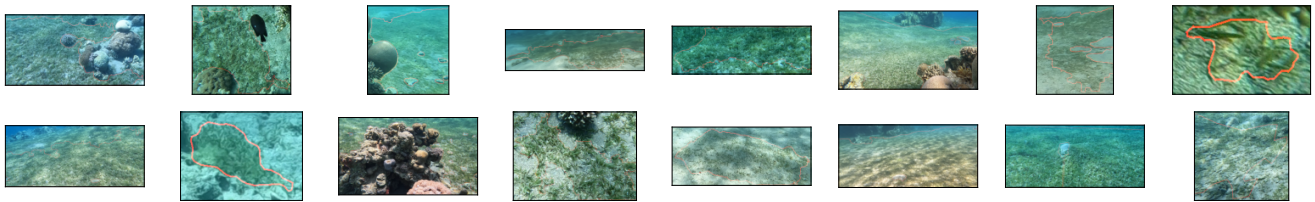
Algae Covered Substrate  Substrate covered in turf algae or other macroalgae, including fleshy algae and Turbinaria.



Hard Substrate  Hard substrate that is part of the reef, which can not be identified into any of the other classes. Includes rocks and heavily decomposed coral skeletons. Also includes human-made structures (underwater infrastructure) such as pipes, pier columns, coral nursery tables, buoys and their lines, as well as boat anchors.



Seagrass 



C. Benchmarking

All models are implemented in PyTorch, with convolutional architectures (UNet++, DeepLabV3+) built using the segmentation-models-pytorch library [37] and transformer-based architectures (SegFormer, DPT) implemented using the Hugging Face transformers library [80]. Parameter efficient fine-tuning with LoRA was done using the Hugging Face PEFT library [52]. Training and evaluation are conducted on a system with a h100 GPU. The code used to train and evaluate the models is available on GitHub. ‡

For the DeepLabV3+ model with a ResNet-50 backbone we follow the training strategy of [21], training for 1000 epochs with a batch size of 16 using stochastic gradient descent (SGD) with an initial learning rate of $1e-3$, momentum of 0.9, weight decay of $1e-4$, and a polynomial learning rate scheduler with a power of 0.9. The model is optimized using a cross-entropy loss. During training, images are randomly scaled within a range of 0.75 and 2, flipped horizontally with a 0.5 probability and randomly cropped to 768×768 pixels. All input images are then normalized using the ImageNet mean and standard deviation. When training on the train+validation set and evaluating on the test set, we trained the model for 700 epochs based on the optimal validation set mIoU.

For the UNet++ model with a ResNet-50 backbone we follow the training strategy of [86], training for 1000 epochs with a batch size of 16 using stochastic gradient descent (SGD) with an initial learning rate of $1e-2$, momentum of 0.99, weight decay of $3e-5$ and polynomial learning rate scheduler with a power of 0.9. The model uses a combination of cross-entropy and Dice loss, with equal weighting. During training, images are randomly cropped to 512×512 pixels and flipped horizontally with a 0.5 probability. Input images are normalized using the ImageNet mean and standard deviation. When training on the train+validation set and evaluating on the test set, we trained the model for 925 epochs based on the optimal validation set mIoU.

For the SegFormer models, training is conducted following [81], using a batch size of 8 when using the MiT-B2 backbone and 4 when using the MiT-B5 backbone for 1000 epochs, using the AdamW optimizer with an initial learning rate of $6e-5$ (multiplied by 10 when using LoRA), weight decay of $1e-2$ and polynomial learning rate scheduler with a power of 1. During training, images are randomly scaled within a range of 1 and 2, flipped horizontally with a 0.5 probability and randomly cropped to 1024×1024 pixels. Input images are normalized using the ImageNet mean and standard deviation. For evaluation, a non-overlapping sliding window strategy is employed, using a window size of 1024×1024 and a stride of 1024. When performing parameter efficient fine-tuning with LoRA, a rank size and alpha of 128 was used. We tried using ranks of 64 and 256, however that had no significant effect on the results. When training on the train+validation set and evaluating on the test set, we trained the model for 265 epochs for the MiT-B2 backbone, 75 epochs for the MiT-B5 backbone and 765 epochs for both backbones trained with LoRA. To increase the augmentation strength, we change the crops of size 1024 to random resized crops of scale 1.02 to 2.0 with aspect ratio 3/4 to 4/3, include random rotations of up to 15 degrees (cropped to exclude non-image areas), as well as random color jitters of contrast, saturation, and brightness between 0.8 and 1.2, and hue changes between -0.05 and 0.05. This model is trained for 100 epochs.

For the Linear DINOv2 model we followed the implementation described in the DINOv2 paper [59]. The model is trained for 1000 epochs with a batch size of 16, using the AdamW optimizer with an initial learning rate of $5e-5$ ($5e-6$ for the backbone layers), weight decay of $1e-2$, and a polynomial learning rate scheduler with a power of 1. During training, images are resized to 1036×518 pixels, randomly horizontally flipped with a 0.5 probability, and normalized using the ImageNet mean and standard deviation. For validation and testing, images are resized to the same dimensions and normalized similarly, while keeping the ground truth masks at their original size. The performance of the model is then evaluated over the interpolated prediction mask at the original dimension of the mask (2048×1024). When training on the train+validation set and evaluating on the test set, we trained the model for 155 epochs based on the optimal validation set mIoU.

For the DPT models [66] with a DINOv2 backbone, we followed the training procedure inspired by the Linear DINOv2 approach described above and the implementation of DepthAnythingV2. As such the model is trained for

‡<https://github.com/eceo-epfl/coralscapesScripts>

1000 epochs with a batch size of 16 when using the base backbone and 8 when using the giant backbone, using the AdamW optimizer with an initial learning rate of $5e-5$ ($5e-6$ for the backbone layers; multiplied by 10 when using LoRA), weight decay of $1e-2$, and a polynomial learning rate scheduler with a power of 1. During training, images are resized to 1036×518 pixels, randomly horizontally flipped with a 0.5 probability, randomly cropped to size 518×518 and normalized using the ImageNet mean and standard deviation. For validation and testing, images are resized to the same dimensions and normalized similarly, while keeping the ground truth masks at their original size, using the same performance evaluation procedure as for the Linear DINOv2 model. When training on the train+validation set and evaluating on the test set, we trained the model for 165 epochs for the DINOv2-Base backbone, 55 epochs for the DINOv2-Giant backbone and 365 epochs for both backbones trained with LoRA.

D. UCSD Mosaics Transfer Learning

All models are trained using the same optimizer and learning rate setting as in the Coralscapes benchmarking experiments. During training, only random horizontal and vertical flips of the 512×512 px patches are used during training as data augmentations. The models are trained for 100 epochs in the 5 & 10 labels per image setting, and for 200 epochs otherwise. Example training labels in each of the settings are shown in Figure 9, and example model outputs for DeepLabV3+ trained in the sparsest setting are shown in Figure 10.

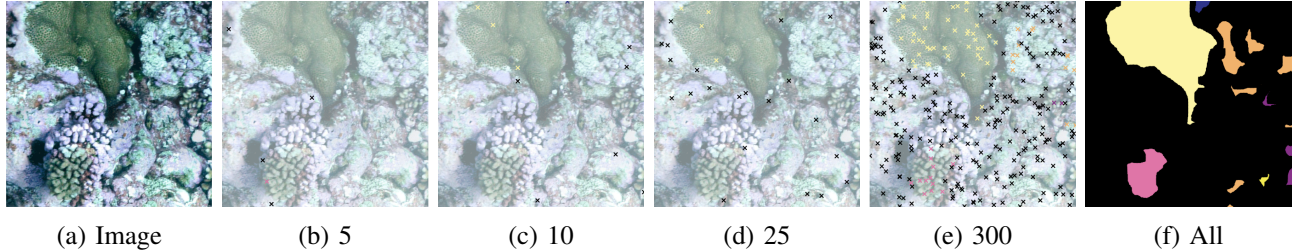


Fig. 9: Example training sample for the UCSD mosaics transfer setting.

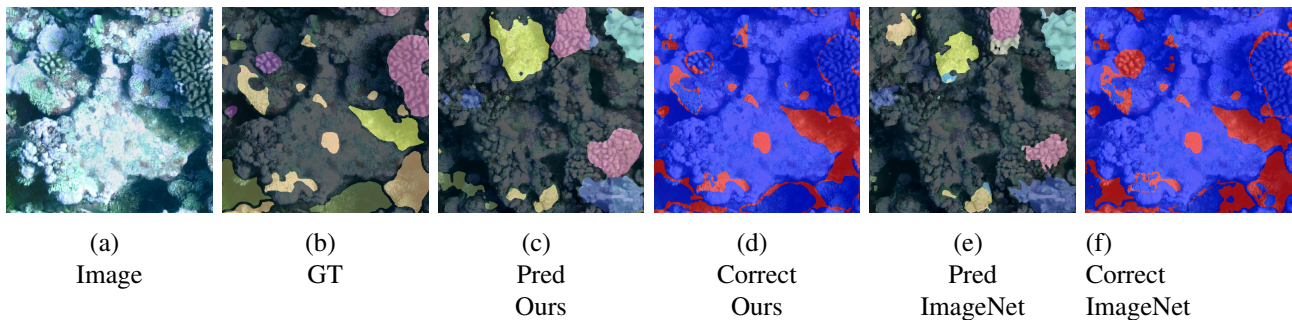


Fig. 10: Example predictions of DeepLabV3+ pre-trained on Coralscapes and an off-the-shelf (ImageNet) model, trained with 5 labels per image.

E. Crown-of-Thorns Starfish Survey

The CSIRO Crown-of-Thorns dataset [50] was originally a challenge on Kaggle. As the competition has finished, the original test dataset is no longer available. The segmentation models that were pre-trained on Coralscapes were saved and the last (output) layer manually reshaped from 39 classes to 2 classes, with the ‘COTS’ neuron being the ‘Crown of Thorn’ neuron from Coralscapes, and the other neuron the mean of the remaining 38 Coralscapes classes. We train the segmentation models for 10 epochs on square patches of size 720, with otherwise the same hyperparameters as training on Coralscapes. At inference, a sliding window with stride 280 is used for prediction. As a baseline, we train YoloV8-L [75] for 100 epochs using the ‘ultralytics’ Python library [40], in the default settings, which initializes from a pre-trained model and applies a range of image augmentations during training.

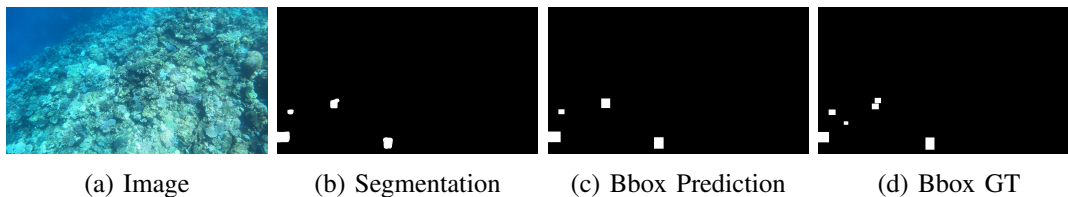


Fig. 11: Example of how the images are transformed to bounding boxes when using SegFormer in the COTS detection task.

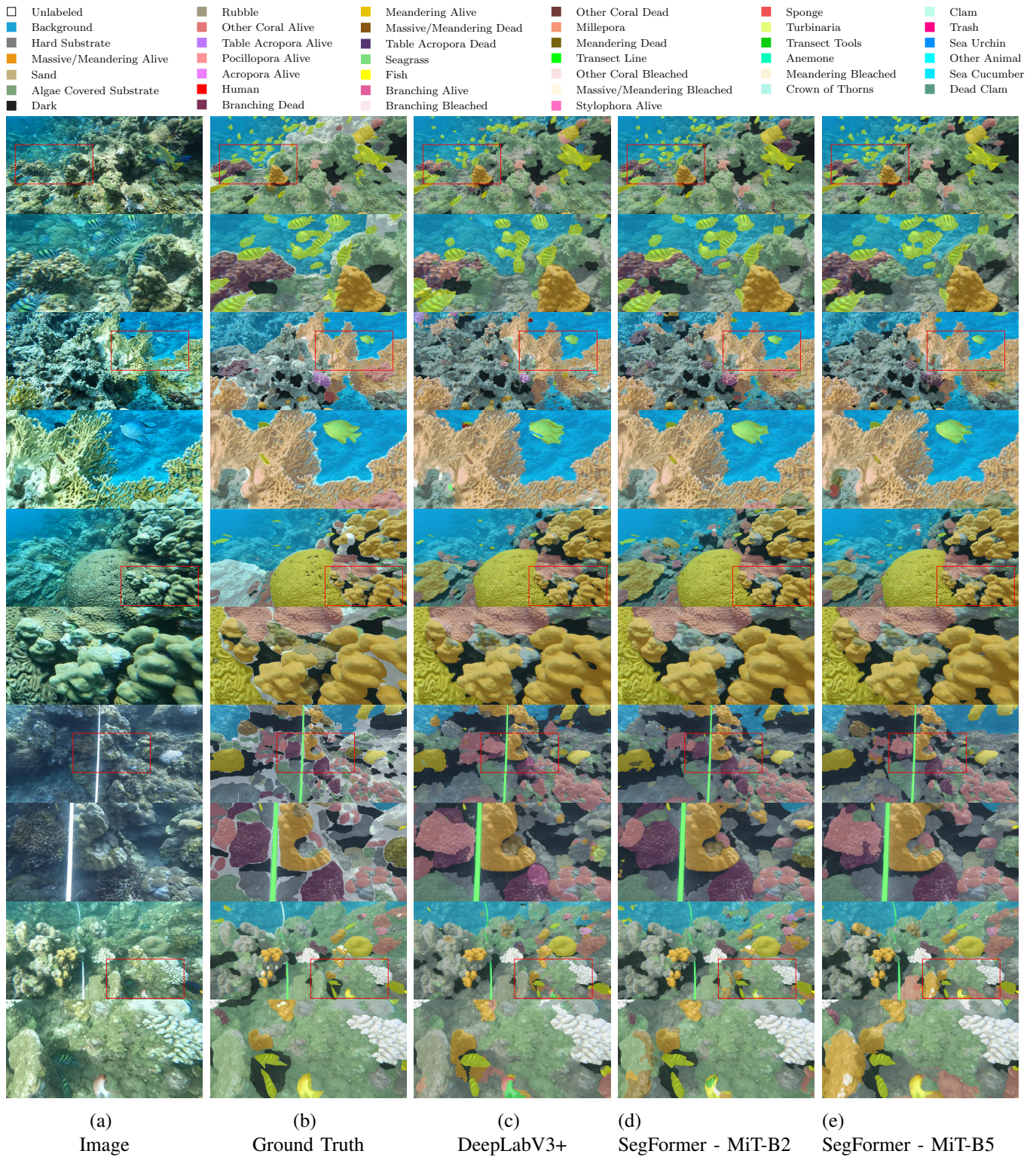


Fig. 12: Additional qualitative samples from the Coralscapes test set.