

Human-in-the-Loop Segmentation of Multi-species Coral Imagery

Scarlett Raine^{1,2} *Graduate Student Member, IEEE*, Ross Marchant³,
Brano Kusy² *Member, IEEE*, Frederic Maire¹, Niko Sünderhauf¹ *Member, IEEE*
and Tobias Fischer¹ *Senior Member, IEEE*

Abstract

Marine surveys by robotic underwater and surface vehicles result in substantial quantities of coral reef imagery, however labeling these images is expensive and time-consuming for domain experts. Point label propagation is a technique that uses existing images labeled with sparse points to create augmented ground truth data, which can be used to train a semantic segmentation model. In this work, we show that recent advances in large foundation models facilitate the creation of augmented ground truth masks using only features extracted by the denoised version of the DINOv2 foundation model and K-Nearest Neighbors (KNN), without any pre-training. For images with extremely sparse labels, we present a labeling method based on human-in-the-loop principles, which greatly enhances annotation efficiency: in the case that there are 5 point labels per image, our human-in-the-loop method outperforms the prior state-of-the-art by 14.2% for pixel accuracy and 19.7% for mIoU; and by 8.9% and 18.3% if there are 10 point labels. When human-in-the-loop labeling is not available, using the denoised DINOv2 features with a KNN still improves on the prior state-of-the-art by 2.7% for pixel accuracy and 5.8% for mIoU (5 grid points). On the semantic segmentation task, we outperform the prior state-of-the-art by 8.8% for pixel accuracy and by 13.5% for mIoU when only 5 point labels are used for point label propagation. Additionally, we perform a comprehensive study into the impacts of the point label placement style and the number of points on the point label propagation quality, and make several recommendations for improving the efficiency of labeling images with points.

¹QUT Centre for Robotics, Australia {sg.raine, f.maire, niko.suenderhauf, tobias.fischer}@qut.edu.au

²CSIRO Data61, Australia {scarlett.raine, brano.kusy}@csiro.au

³Image Analytics, Australia ross.g.marchant@gmail.com

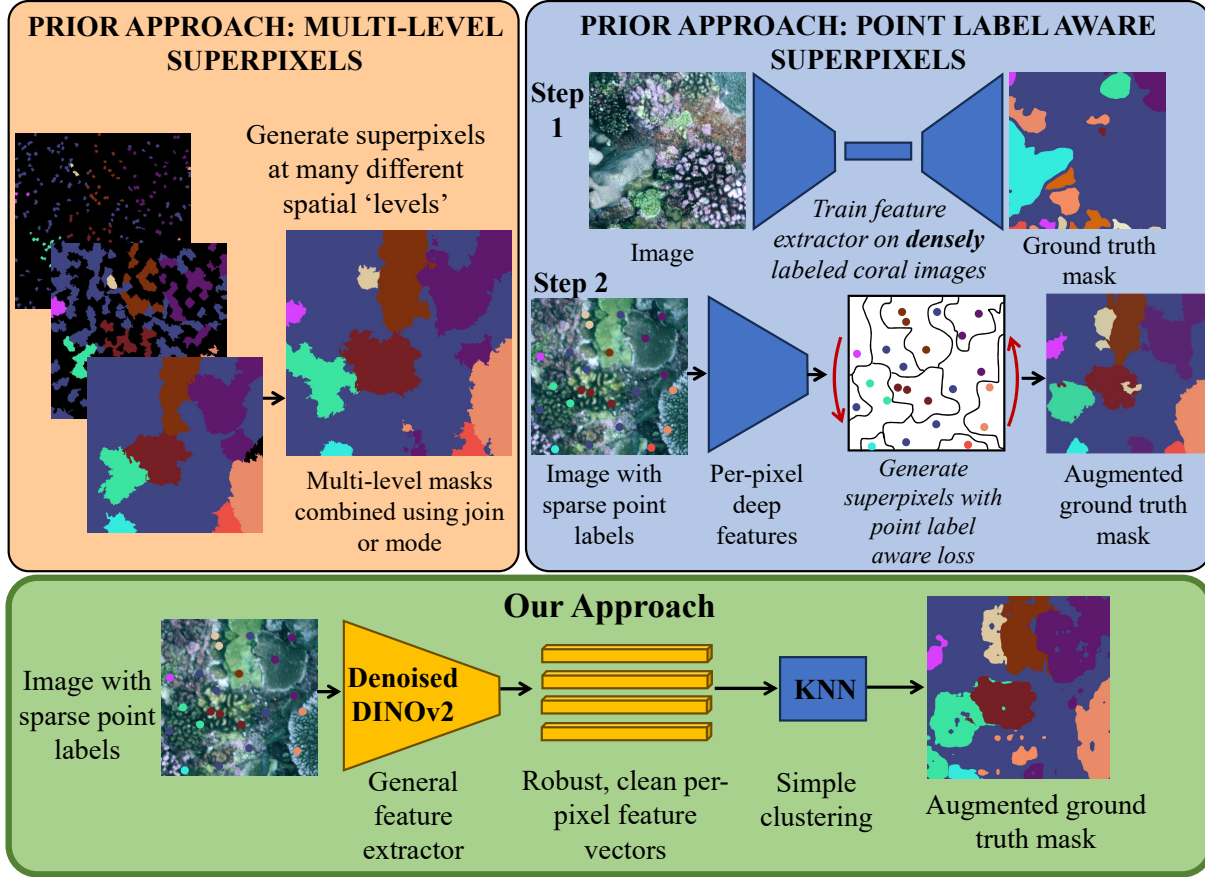


Fig. 1. The proposed point label propagation technique utilizes the DINOv2 foundation model without any fine-tuning to create augmented ground truth masks for intricate coral images. **Top left:** Previous methods depended on layering superpixels that contained point labels [1]–[3]. **Top right:** A more recent approach involved pre-training a CNN feature extractor on densely labeled coral images and propagating point labels with a custom, point label aware superpixel function [4]. **Bottom:** In contrast, our approach employs KNN to group features derived from the denoised version of the DINOv2 foundation model without further training on coral imagery.

Index Terms

Environmental Monitoring, Scene Understanding, Robotics, Human-in-the-Loop, Foundation Models, Semantic Segmentation

I. INTRODUCTION

INFORMED marine ecosystem management relies on data gathered at various spatial and temporal scales [5]. Marine surveys are increasingly being conducted using autonomous underwater and surface vehicles [6]–[8]. However, these methods produce vast amounts of

seafloor images that need to be analyzed to produce meaningful statistics such as coverage estimates of various substrates and coral species [9]–[11].

One method for automated image analysis is semantic segmentation, a computer vision task that predicts the class of every pixel in a query image [12]. In recent years, it has become common practice to train deep learning models for performing semantic segmentation. These models are often trained using full supervision, *i.e.* on pairs of images and dense ground truth masks, where the mask is comprised of a class label for every corresponding pixel in the paired image [13]. These dense masks are typically created by human experts who manually label each pixel.

However, underwater imagery has visually unique characteristics which complicate the process of labeling images. Coral images frequently exhibit high complexity, and feature unclear boundaries, significant color and texture variation among species, and low clarity [14]–[17]. The complex and intricate nature of coral images and the challenges in correctly classifying coral species necessitate annotation of images by marine scientists, thus hindering use of standard computer vision annotation tools like the crowd-sourcing provider Amazon Turk¹.

Historically, marine scientists have labeled underwater images using the Coral Point Count (CPC) [18] methodology, in which a scientist specifies the class of a set quantity of randomly distributed or grid-spaced sparse pixels in each image. We call these annotated pixels *point labels* [4]. The CPC method typically involves labeling 100-300 point labels per image. The extremely sparse label setting that we target in this paper is critical for ecologists who often have limited budget and time constraints for monitoring projects and are unable to label 100-300 points per image [19]. The sparse label setting can also occur during field expeditions, when ecologists might be performing surveys in a new location, for a specific task or to detect certain species. In this case, it is necessary to quickly train, re-train or fine-tune models to perform semantic segmentation, often overnight or between transect surveys, but it is prohibitively time-consuming and costly for domain experts to label every pixel or even 300 pixels in each image [19]. By contrast, it could be feasible for experts to quickly label 5 points per image and use a point label propagation algorithm like ours to create dense augmented ground truth for training or fine-tuning a semantic segmentation model.

Although there are large quantities of CPC data available in point grid and random formats

¹<https://www.mturk.com/>

[20], [21], the most efficient way of labeling points for training deep learning models to perform semantic segmentation has not been investigated, and could have significant impacts for annotation costs.

Recently, superpixel algorithms based on RGB color values [1]–[3] and deep features [4] have been leveraged for propagating point labels into dense, pixel-wise augmented ground truth masks. These masks are then used as the supervisory signal when training deep neural networks to perform semantic segmentation of underwater images. Raine *et al.* [4] presented an innovative point label aware approach to superpixels, clustering pixels into segments using the per-pixel deep CNN features and the RGB values. Although this algorithm advanced the state-of-the-art, the deep features are extracted by a CNN trained on coral imagery, and experienced performance issues when only a small number of point labels were available.

In this study, we address the scenario where only a very small number of labels are available. This situation is crucial because marine survey projects frequently have constrained budgets for data labeling [5]. Moreover, processing survey data often requires swiftly fine-tuning and retraining models in the field to adapt to new locations, species or environmental conditions [5]. In this work, we propose using the DINOv2 foundation model [22], [23] to generate per-pixel embeddings. We then employ the straightforward K-Nearest Neighbor algorithm to create the augmented ground truth, surpassing the state-of-the-art with a limited number of point labels. Additionally, we show further performance enhancements by implementing a human-in-the-loop point selection strategy, utilizing the expertise of human experts to reduce uncertainty in the embedding space of the KNN.

This work establishes the utility of general foundation models for multi-species segmentation of domain-specific underwater imagery, with particular focus on label efficiency for the extremely sparse label setting (Fig. 1). Our contributions are summarized as follows:

- 1) We propose using a general-purpose foundation model to produce per-pixel deep features for coral images, demonstrating that these features are discriminative without training or fine-tuning on coral imagery. When combined with the basic K-Nearest Neighbors algorithm, these features are sufficient for creating accurate augmented ground truth masks, eliminating the need for custom superpixel algorithms.
- 2) We consider the extremely sparse point label setting, *i.e.* 5-25 points per image, and present a novel human-in-the-loop labeling regime, which integrates the expertise of the marine scientist with the model’s introspective uncertainty to identify informative locations for

point labels. Our approach outperforms the prior state-of-the-art method by 14.2% pixel accuracy and 19.7% mIoU for the 5 point label setting, and by 8.9% and 18.3% when there are 10 points labeled in each image.

- 3) If human-in-the-loop labeling is not available, we find that leveraging the DINOv2 denoised features with a KNN results in improvements over the state-of-the-art for propagation of small quantities of point labels per image. We see pixel accuracy improve by 2.7% and an improvement of 5.8% for mIoU (5 point labels per image); and 2.3% in pixel accuracy and 10.0% in mIoU (10 point labels per image) on UCSD Mosaics.
- 4) These improvements in point label propagation are reflected in the semantic segmentation task. A DeepLabv3+ model trained on augmented ground truth masks generated using DINOv2, KNN and our human-in-the-loop labeling regime significantly outperforms the prior works: pixel accuracy improves by 8.8% and mIoU improves by 13.5% (5 point labels per UCSD Mosaics image).
- 5) We conduct comprehensive experiments to assess the impact of the quantity and placement of point labels on the point propagation task, offering valuable recommendations for efficient annotation.

To foster future research in this area, we make our code publicly available².

II. RELATED WORK

Performing marine surveys with autonomous surface and underwater vehicles enables the collection of significant quantities of images [6], [7], [24], [25]. Automating the analysis of underwater imagery requires innovative approaches based on deep learning, robotics, computer vision and specialized expertise in underwater ecosystems [10], [26]. This section reviews and analyzes prior methods for semantic segmentation of underwater imagery and point label propagation (Section II-A), advances in large foundation models (Section II-B), and human-in-the-loop fundamentals (Section II-C).

A. Semantic Segmentation of Coral Imagery

Underwater semantic segmentation is challenging due to the dynamic, unrestricted environment of the ocean [17]. Typically images are deteriorated by noise, turbidity, scattering and attenuation

²<https://github.com/sgraine/HIL-coral-segmentation>

of sunlight, low-light conditions, blur, and changes in coloration due to depth [16], [17]. For example, image characteristics vary significantly both within and between datasets due to lighting (whether artificial lighting or natural lighting is employed), camera settings used, and post-processing of images. Coral species often exhibit strong visual similarity between different species, and are typically intricate and highly textured [15], [19], [20]. Underwater image characteristics, along with the lack of distinct “objectness” for overlapping corals, make the underwater semantic segmentation task a uniquely challenging problem.

Many works which perform fully supervised segmentation of underwater coral imagery train the model with images and their corresponding densely labeled, pixel-wise ground truth masks [27]–[30], [30]–[33]. Some approaches have proposed interactive labeling tools which use deep learning to assist annotation, but these approaches do not target label efficiency or propose informative label locations [30], [34], [35]. The Machine learning Assisted Image Annotation (MAIA) method proposes novel regions for labeling, but relies on an autoencoder pre-trained on coral images [36]. TagLab is a data labeling tool designed to speed up the pixel-wise annotation of large orthoimages³, however it is based on a model trained using 15,000 coral images with dense ground truth labels [38]. Recently, a human-in-the-loop iterative labeling approach was proposed which combines expert and crowd-sourced non-expert annotations, however their method is only for bounding box detection of marine species [39].

There are limited algorithms in the literature for weakly supervised segmentation of corals [1]–[4], [40], [41]. These rely on custom-designed superpixel algorithms which propagate sparse point labels to obtain corresponding pixel-wise ground truth. The multi-level superpixel method generates superpixel segments informed by the RGB features. This algorithm repeatedly generates the segments at a variety of spatial scales, each time labeling each segment as the class of the point label inside before combining the various “levels” of superpixels into a single augmented ground truth [1]–[3]. More recently, Point Label Aware Superpixels [4] was proposed and introduced a novel superpixel algorithm which clusters pixels with deep features and creates superpixel segments informed directly by the point labels.

Previous superpixel methods depend on the availability of a sufficient number of points and suffer from degraded performance in very sparse settings. There is potential to utilize the recent

³Orthoimages or orthomosaics are geometrically corrected images which have been adjusted to remove distortions due to the camera, water and seafloor topography. Orthoimages can be generated from a series of overlapping images taken by an underwater camera or robotic vehicle and then joined together into a mosaic to cover a larger area [37].

innovations in large foundation models for a more general, simplified approach to point label propagation.

B. Large Foundation Models

Recently there has been considerable research effort towards developing large foundation models which learn robust, discriminative feature embeddings that are task-agnostic [22], [42], [43]. Foundation models are trained on extensive datasets and designed to acquire highly generalized representations, enabling them to apply their knowledge to tasks and data beyond the training scope [42].

One such foundation model, DINOv2 [22], was trained on 142 million images, and is based on a Vision Transformer (ViT) model [44] trained with a discriminative self-supervised loss function [22]. This loss function combines an image-level objective function for features extracted from a student-teacher framework and a patch-level objective, in which patches from the input image are masked and the model must predict the missing image regions [22].

Some studies focus on tailoring foundation models to specific tasks or contexts, such as for recognition of a user’s pet [45], or by training a decoder or adapter on the desired task [46]. In the context of plant phenotyping, modified foundation models were evaluated for instance segmentation, leaf counting and disease classification, however the methods designed specifically for these tasks outperformed the modified foundation models [46]. In medical image analysis, [47] show that DINOv2 has cross-task generalizability and report competitive results when its features are applied with KNN for disease classification. Other research has explored self-supervised object localization without using labels [48], but these methods do not achieve segmentation of the entire image.

While some works have explored the utility of foundation models for specialized problems [46], [47], [49], using DINOv2 for the task of underwater coral segmentation has not been investigated. As highlighted by Section II-A, coral imagery has unique and challenging visual characteristics, including detailed textures, poorly defined boundaries, and overlapping species, as seen in Fig. 2 [4], [20]. When applied to common imagery, such as on an example image from the Cityscapes [50] dataset, the DINOv2 model is able to effectively segment image, however when used directly on the coral images, the model is not able to generate useful segments (Fig. 2). This suggests that out-of-the-box usage of DINOv2 on this domain-specific task is not effective. However, our work investigates whether the feature representation of the foundation

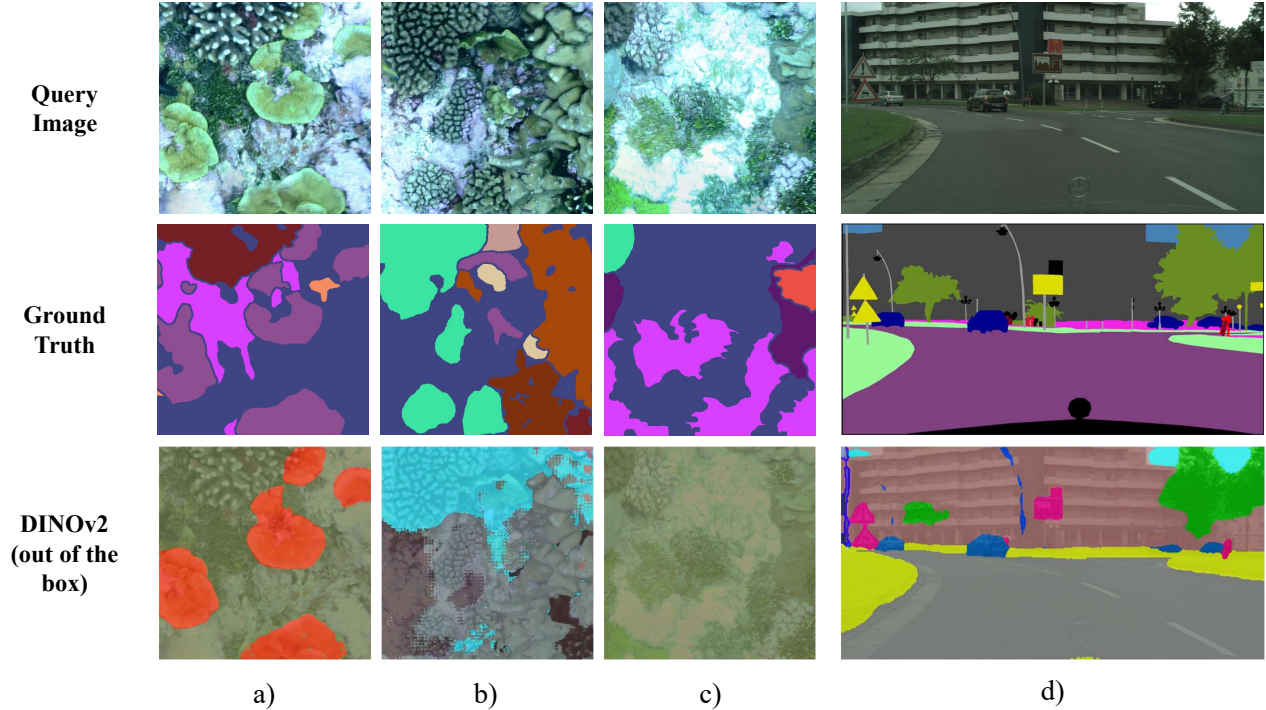


Fig. 2. A comparison of segmentation by DINOv2 [22] on UCSD Mosaics [2], [51] coral images (a, b, and c) and an image from the Cityscapes [50] dataset (d). Note that in example c), the DINOv2 model does not produce any segments. This demonstrates that DINOv2 cannot be directly used in domain-specific applications. This work instead proposes leveraging the deep feature space in combination with a nearest neighbor classifier to perform point label propagation.

model DINOv2, trained on general images, can generate meaningful feature embeddings for coral images. We propose leveraging the deep embedding space of DINOv2 alongside sparse domain-expert point labels for point label propagation in coral imagery.

C. Human-in-the-Loop

Human-in-the-Loop refers to machine learning algorithms and pipelines that enable humans to directly provide feedback to and engage with a model [52]. More specifically, the Human-in-the-Loop subfield of Interactive Machine Learning outlines a framework where control is shared between humans and models; humans provide input to the system in a direct, frequent and interactive manner [52], [53].

Although models have previously been used to predict labels on unseen data in the field of ecology [54], [55], the application of foundation models within an interactive labeling framework for coral point label propagation has not yet been explored.

To our knowledge, there is no existing approach in the literature for semantic segmentation of multi-species coral imagery that integrates the power of task-agnostic foundation models with domain-specific, annotation-efficient labeling. This presents an opportunity to reduce the time and costs associated with manual annotation of complex, domain-specific imagery, while improving point label propagation when there are extremely limited labels available.

III. METHOD

A. Method Overview

This section provides an overview of our proposed point label propagation approach. We leverage the denoised version of the DINOv2 foundation model by [23], based on [22]. We cluster pixels in the deep embedding space with K-Nearest Neighbors, yielding the augmented ground truth mask.

Our method expects a photo-quadrat coral image and a collection of sparse point labels as input, producing a dense augmented ground truth mask as output. The input point labels may be randomly distributed across the image, arranged evenly in a grid⁴, or designated by our proposed human-in-the-loop framework (Section III-B).

From our image, we generate a set of feature vectors, each 768 in length, extracted from the denoised DINOv2 feature extractor [23]. This extractor produces a deep feature for every 14x14 pixel patch in the input image. Additionally, it generates a ‘CLS’ token for the entire image, which is not used in our method. We then perform spatial upsampling on the feature vectors using bilinear interpolation to create a deep feature vector for each pixel in the input image, after which we L2 normalize the per-pixel feature vectors.

We obtain a set of sparse labeled pixels L and retain the normalized feature embeddings $\{\mathbf{v}_1, \dots, \mathbf{v}_l, \dots, \mathbf{v}_L\}$ for these pixels. In addition, we store $X = \{(x_1, y_1), \dots, (x_L, y_L)\}$ where (x_l, y_l) are the pixel coordinates of \mathbf{v}_l . We find the cosine similarity of the feature embedding \mathbf{v}_l for $l \in \{1, \dots, L\}$ and the feature embedding \mathbf{v}_p for every other pixel $p \neq l$ in the image:

$$\text{sim}(\mathbf{v}_p, \mathbf{v}_l) = \mathbf{v}_p \cdot \mathbf{v}_l. \quad (1)$$

⁴If the grid cannot evenly accommodate the total number of points, the nearest possible arrangement is used; for example, in the case of 5 point labels, the point labels are spaced in a 2x2 grid with a single point in the center. For the 10 point label case, we use 3 rows of points, where the first and third rows contain 3 point labels and the middle row has 4 evenly spaced points.

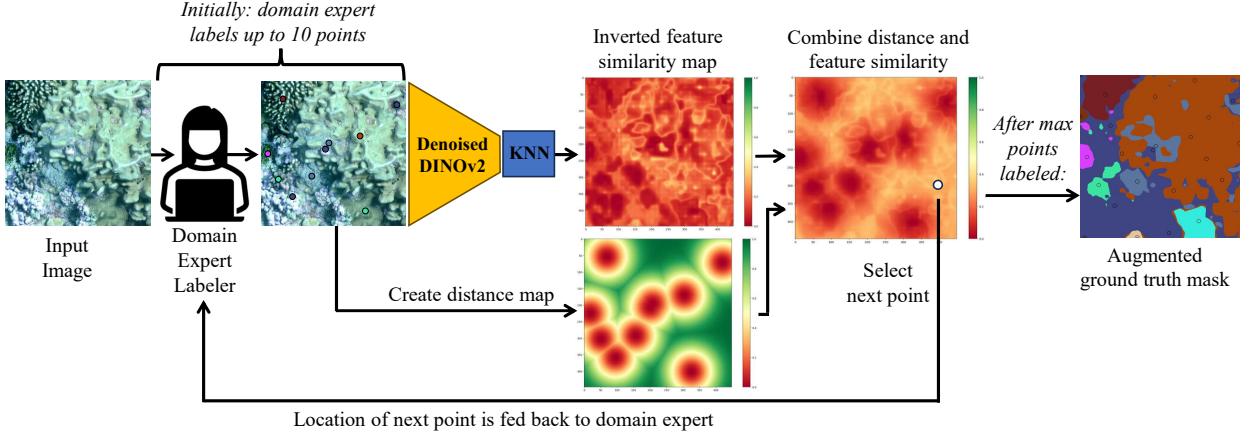


Fig. 3. Schematic of Proposed Human-in-the-Loop Labeling Approach. We combine domain expert knowledge with the model’s internal uncertainty to improve point label selection. The process starts with inputting a coral image and having a marine scientist label up to 10 points centrally located on the largest instances (refer to Section V-B for an analysis of this value). A feature similarity map is then generated by computing cosine similarities between the labeled points and all other pixels. To promote exploration, we use a distance map in conjunction with the similarity map to create a combined probability mask for pixel selection. The chosen pixel is then sent back to the marine scientist for labeling, and the KNN is updated accordingly. Once the maximum number of points has been labeled, an augmented ground truth mask is created for use in training a semantic segmentation model.

We obtain a dense ground truth mask by performing K-Nearest Neighbors with $k = 1$. Note that although we evaluated various values for k this did not lead to any improvement for $k > 1$, as discussed in Section V-B and in Fig. 10.

B. Human-in-the-Loop Pixel Proposal

To enhance our propagation accuracy in cases with extremely sparse data, we design a novel labeling regime (Fig. 3). Unlike previous methods that have utilized random or grid-spaced sparse point labels, we treat the point labeling task as a human-in-the-loop process. We assume the availability of a marine scientist who can collaboratively label a specific number of points, which will then be incorporated into our DINOv2 and KNN point label propagation approach.

To identify informative points for the marine scientist to annotate, we analyze the cosine similarity between the features of labeled and unlabeled pixels in the DINOv2 deep embedding space. Initially, we ask the marine scientist to label up to 10 pixels located centrally in the largest instances visible in the image (see Section V-B for an ablation study on this number). For instances requiring more than 10 labeled points, the human-in-the-loop regime iteratively

proposes one point at a time for labeling, focusing on the areas of the image with the highest uncertainty. This uncertainty is quantified as the cosine similarity to the nearest labeled pixel.

To implement this, we apply the method described in Section III-A to obtain, upsample, and normalize the per-pixel feature embeddings. We then calculate a cosine similarity map (Eq. 1) between the initially labeled pixels and every other pixel in the image. By inverting this map, we assign a higher probability of selection to pixel locations that exhibit low cosine similarity to the nearest labeled pixel:

$$C(x, y) = 1 - \max_{l \in \{1, \dots, L\}} \text{sim}(\mathbf{v}_q, \mathbf{v}_l), \quad (2)$$

where \mathbf{v}_q is the feature vector at location (x, y) .

To promote exploring the entire image, we generate a probabilistic distance map based on the labeled pixels. We calculate the Euclidean distance transform on a binary mask that indicates the positions of the labeled pixels, where the initial count is $L = 10$:

$$D(x, y) = \min_{(x', y') \in X} \sqrt{(x - x')^2 + (y - y')^2}, \quad (3)$$

where $X = \{(x_1, y_1), \dots, (x_L, y_L)\}$, which stores the pixel coordinates of the labeled points.

We then perform Gaussian smoothing over the distance transform and tune the smoothing parameter σ in the ablation study in Section V-B:

$$D_{\text{smooth}}(x, y) = 1 - \exp\left(-\frac{D(x, y)^2}{2\sigma^2}\right). \quad (4)$$

We combine the probabilistic cosine similarity map with the distance map, and weight the two terms with λ (see hyperparameter tuning in Section V-B):

$$M(x, y) = \frac{D_{\text{smooth}}(x, y) + \lambda C(x, y)}{\lambda + 1}. \quad (5)$$

From the combined map we identify the next pixel for annotation by the expert by selecting the location $(\hat{x}, \hat{y}) = \arg \max_{(x, y)} M(x, y)$ corresponding to the highest selection probability in M .

C. Semantic Segmentation

After point label propagation (Stage One), we assess the augmented ground truth masks by training a fully supervised model for semantic segmentation (Stage Two). The model chosen for

Stage Two can be tailored to meet the computational and inference time needs for deployment. In this study, we select the DeepLabv3+ model [56], a widely used architecture for semantic segmentation, and train it using the augmented ground truth masks as outlined in Section IV-A.

IV. EXPERIMENTAL SETUP

In this section, we outline the details of our implementation in Section IV-A, describe the evaluation datasets in Section IV-B, and define the evaluation metrics in Section IV-C.

A. Implementation

1) *Stage One: Point Label Propagation:* All experiments in this work are completed on a Quadro RTX 6000 GPU, and we calculate point label propagation times in Table I with respect to this GPU. We implement our presented approach with Python and PyTorch [57]. In addition, we employ the Faiss library to enable faster implementation of K-Nearest Neighbors on GPU [58]. The denoised DINOv2 model and implementation is from [23].

2) *Stage Two: Semantic Segmentation:* We evaluate the propagated ground truth masks by training a model for semantic segmentation. We use TensorFlow to train the DeepLabv3+ model [56] and apply data augmentation techniques, including random horizontal and vertical flipping, and adjust the gain and gamma values within the range of 0.8 – 1.2. The training process spans 500 epochs using the Adam optimizer with a learning rate of 0.001. We report results on the test dataset based on the best epoch.

Unlike prior works [2] and [4], for Stage Two we include all classes for training and evaluation. Previously, prior works have excluded the “unknown” class in the UCSD Mosaics dataset (Section IV-B) from training and test. In this work, due to the sparsity of the labels considered in this setting, we include this class during training and test to ensure a fair evaluation of all models. If the “unknown” class is excluded, this results in model under-fitting and misrepresentation of the model performance.

B. Datasets

In this work, we use the UCSD Mosaics dataset, which contains multi-species coral images annotated with pixel-wise ground truth masks [2], [51]. In [2], the authors provide a version of the dataset where the large orthomosaics have been saved as smaller images, however we have identified some corrupted ground truth masks and have excluded them (this resulted in 219 being

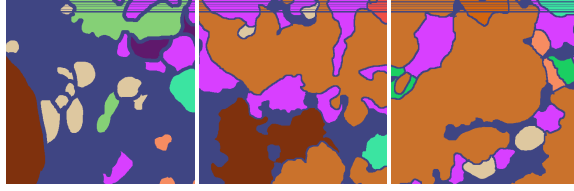


Fig. 4. Some of the ground truth masks in the UCSD Mosaics dataset showed signs of corruption, as illustrated at the top of these examples. The corresponding images and ground truth masks were removed from the training and test datasets.

removed from the training set, resulting in 3,974 images, and 32 were removed from the test set, leaving 696 images; examples of the corruption can be seen in Fig. 4). We include a list of the corrupted filenames alongside our publicly available code⁵. Each image measures 512 by 512 pixels and features 33 types of corals along with an class called ‘unknown’ or ‘unlabeled’. To maintain consistency with [2] and [4], we disregard this class during evaluation of the point label propagation task. We imitate the role of the domain expert in our human-in-the-loop labeling framework by obtaining the ground truth point label corresponding to the proposed location. For the first 10 pixels, we select the central pixel of the largest instances in the mask.

C. Evaluation Metrics

We employ three commonly used metrics [2], [4], [59] to evaluate the performance of our method against previous approaches:

- 1) Pixel Accuracy (PA), which measures the proportion of correctly predicted pixels out of all predicted pixels,
- 2) Mean Pixel Accuracy (mPA), the average pixel accuracy across all classes, and
- 3) Mean Intersection over Union (mIoU), representing the average of the IoU scores for each class.

For all these metrics, a higher value signifies higher performance. These metrics are used to quantify the performance of both Stage One and Stage Two.

V. RESULTS

We compare our proposed approach to the point label propagation state-of-the-art in Section V-A and then Section V-B provides various ablation studies.

⁵<https://github.com/sgraine/HIL-coral-segmentation>

TABLE I
PERFORMANCE OF STAGE ONE: POINT LABEL PROPAGATION APPROACHES (REFER TO SECTION IV-C FOR METRIC DEFINITIONS), FOR 5 / 10 / 25 / 300 POINT LABELS. ‘F-MSS’ IS FAST MSS [3], ‘PLAS’ IS POINT LABEL AWARE SUPERPIXELS [4], AND ‘D+NN’ IS KNN WITH DENOISED DINOv2 [23] (OURS).

Method	Label Style	PA	mPA	mIoU	Time per Image (s)
		5 / 10 / 25 / 300	5 / 10 / 25 / 300	5 / 10 / 25 / 300	5 / 10 / 25 / 300
F-MSS	Rand.	7.29 / 13.49 / 30.09 / 86.81	6.60 / 12.34 / 29.26 / 82.70	6.55 / 12.11 / 28.53 / 80.12	2.14 / 2.19 / 2.21 / 2.76
F-MSS	Grid	7.94 / 15.58 / 39.18 / 89.98	7.54 / 14.96 / 36.72 / 88.17	7.50 / 14.74 / 35.51 / 86.44	2.43 / 2.45 / 2.36 / 2.96
PLAS - <i>Single</i>	Rand.	48.45 / 55.26 / 65.16 / 86.68	32.03 / 41.44 / 57.65 / 81.74	23.86 / 32.22 / 47.76 / 77.56	1.71 / 2.00 / 2.17 / 1.93
PLAS - <i>Single</i>	Grid	52.08 / 59.09 / 72.96 / 89.28	39.89 / 44.88 / 64.91 / 86.16	30.32 / 36.22 / 58.00 / 82.73	1.55 / 1.80 / 2.06 / 1.81
PLAS - <i>Ens.</i>	Rand.	52.73 / 62.00 / 71.11 / 92.47	36.48 / 49.04 / 63.21 / 89.93	25.91 / 35.6 / 50.46 / 85.45	4.27 / 4.55 / 5.02 / 5.35
PLAS - <i>Ens.</i>	Grid	57.41 / 67.45 / 76.31 / 94.60	44.36 / 54.44 / 69.13 / 92.49	32.89 / 41.19 / 59.82 / 89.38	4.06 / 4.25 / 5.15 / 5.28
D+NN (Ours)	Rand.	55.72 / 64.51 / 75.07 / 88.77	39.94 / 50.91 / 65.80 / 83.84	32.09 / 42.79 / 58.04 / 81.75	4.88 / 4.55 / 4.74 / 4.90
D+NN (Ours)	Grid	60.08 / 69.79 / 78.74 / 89.86	47.85 / 58.39 / 70.05 / 87.41	38.72 / 51.20 / 64.40 / 85.77	4.78 / 4.70 / 4.79 / 4.69
D+NN (Ours)	HIL	71.56 / 76.38 / 81.27 / 89.61	61.46 / 69.87 / 75.91 / 86.45	52.60 / 59.48 / 67.97 / 85.00	4.74 / 4.98 / 20.0 / 273.08

A. Comparison to State-of-the-art Methods

1) *Stage One: Point Label Propagation:* In this section, we evaluate the performance of our approach and compare against a number of state-of-the-art methods, including: Fast Multi-level Superpixel Segmentation (*Fast MSS*) [3], a faster re-implementation of CoralSeg [2], and Point Label Aware Superpixels [4]. In the case of [4], we compare against the single method (*Single*) and also the ensemble of three Point Label Aware algorithms (*Ensemble*).

As demonstrated in Table I, using K-Nearest Neighbors with features extracted by the denoised DINOv2 foundation model [23] for point label propagation surpasses previous methods when dealing with a small number of point labels (5, 10, and 25 per image). When using our human-in-the-loop labeling framework with five point labels, the mIoU increases by 46.1% and 22.6% compared to the Fast MSS (F-MSS) and Point Label Aware Superpixel (PLAS) algorithms, respectively (Fig. 5). Additionally, we see an improvement of 64.3% and 17.3% in pixel accuracy (Fig. 5). Without the human-in-the-loop labeling framework, our method still exceeds the state-of-the-art PLAS by 3.5% in pixel accuracy and 5.7% in mIoU (if 5 grid points are used).

In a scenario not targeted by this paper, where up to 300 point labels are available, our approach shows performance comparable to single classifier methods. However, the ensemble of three Point Label Aware Superpixel classifiers outperforms our approach (Table I).

Our approach, which combines DINOv2 with our human-in-the-loop labeling regime, shows

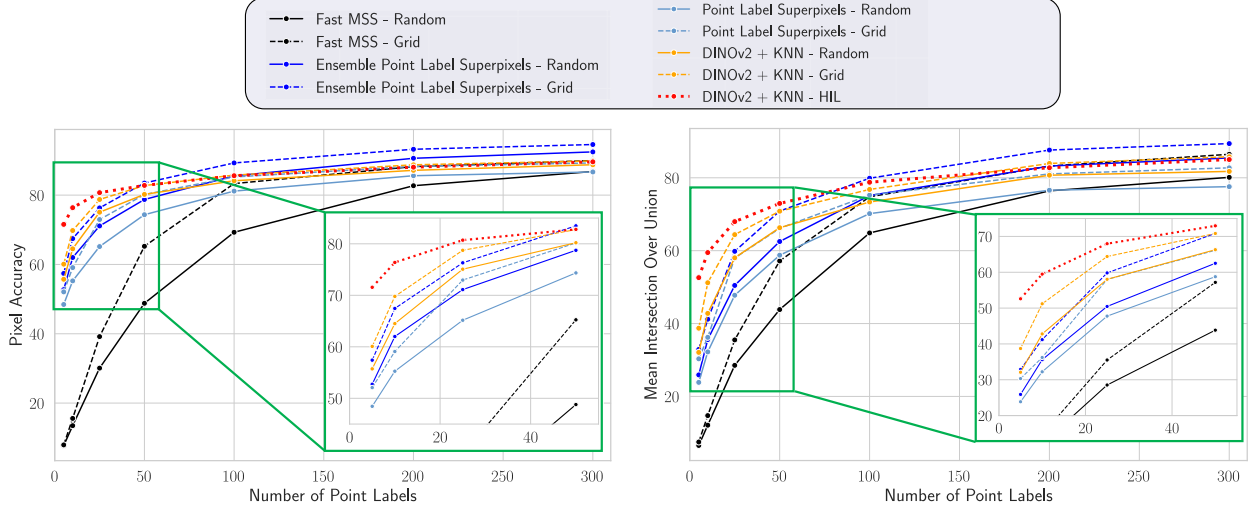


Fig. 5. Point Label Propagation Performance (Pixel Accuracy and Mean IoU): In the graphs, our proposed DINOv2 and KNN method is represented in orange for both random and grid labeling, while the red line shows the performance of the DINOv2 and KNN with the Human-in-the-Loop labeling approach. Our method notably surpasses previous approaches when only a small number of point labels are available, specifically between 5 and 25 points. However, when the number of points increases to 300, the performance of all approaches tends to converge.

comparable computation times per image to the Point Label Aware Superpixel ensemble for small numbers of point labels (Table I). However, when applied to a large number of points *i.e.* 300 points, the computation time becomes excessively high due to the clustering in the deep feature space required to generate the feature similarity map at each iteration (Eq. 2). If there are a large number of point labels available, it is more suitable to use the grid-spaced version. We emphasize that the primary use case for human-in-the-loop labeling is to enhance performance in scenarios with very sparse point labels (5-25 points).

Fig. 6 showcases outputs from our proposed method, illustrating it produces pixel-wise augmented ground truth masks which largely agree with the provided ground truth, even if there are very few sparse point labels available. This figure emphasizes that grid-based sparse labels offer better coverage compared to randomly placed sparse labels. For instance, in row 6, the Fast MSS approach [3] misses an entire beige segment with randomly placed points, while the grid points capture it.

The failure modes of Fast MSS and the Point Label Aware Superpixel approach in the case of 5 pixel labels are also evident in Fig. 6. Fast MSS does not generate useful segments because it only includes segments with point labels in the augmented ground truth mask. When segments lack

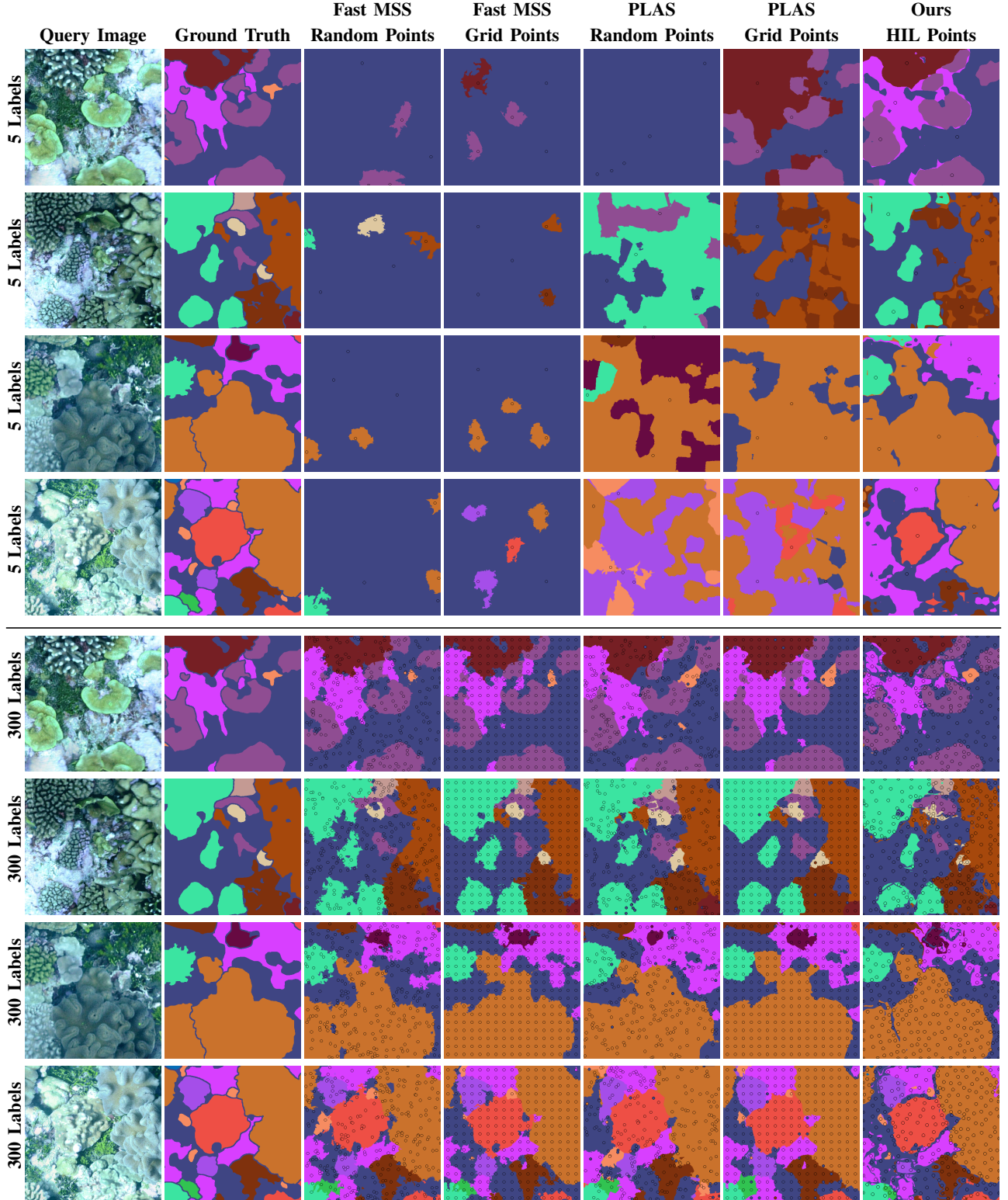


Fig. 6. A comparison is made between Fast MSS [3], Point Label Aware Superpixel (PLAS) [4], and our approach, which combines denoised DINOv2 features [23], K-Nearest Neighbors, and our Human-in-the-Loop labeling scheme. All methods are demonstrated on the same four examples. The top section shows point label propagation with 5 labels, while the bottom section displays propagation with 300 labels. Black circles indicate the point labels used for propagation within each output augmented mask. While all methods effectively propagate ground truth with 300 labels, the Fast MSS and PLAS methods struggle to produce usable augmented ground truth with only 5 labels. In contrast, our approach generates augmented ground truth that accurately resembles the ground truth even with the sparse 5 label setting.

any points (a common occurrence in this setting), they are assigned to the unknown/unlabeled class, resulting in the majority of the mask being this class. Conversely, the PLAS approach labels any superpixel segment without a point based on feature similarity with labeled segments, leading to significant over-prediction of classes. Additionally, PLAS requires sufficient points for the conflict loss function to ensure superpixel boundaries conform accurately to species [4].

Our DINOv2 and KNN approach accurately generates augmented ground truth masks, even with extremely sparse labels. However, a limitation is that small species segments can be missed when there are very few point labels available. For example, in row 1 of Fig. 6, the orange segment is not included in the augmented ground truth. Future work could address this by incorporating mechanisms that emphasize smaller species segments and prevent the model from being biased towards larger instances.

2) *Stage Two: Semantic Segmentation:* We trained a DeepLabv3+ model for semantic segmentation using the training regime and hyperparameters in Section IV-A, with results shown in Table II. The performance improvements for point label propagation are reflected when training a model on augmented ground truth masks. When trained on augmented ground truth masks generated using our DINOv2 and KNN approach with our human-in-the-loop labeling regime, the DeepLabv3+ model outperforms the previous state-of-the-art Point Label Aware Superpixels [4] by 8.8% in pixel accuracy and by 13.5% in mIoU with 5 point labels. Even without the human-in-the-loop labeling regime, our denoised DINOv2 and KNN approach still surpasses prior methods by 6.5% in pixel accuracy and by 6.2% in mIoU (5 point labels). Notably, the DeepLabv3+ model trained on our DINOv2 and KNN approach with grid-spaced labels slightly outperforms prior approaches in the 300 label case, despite not outperforming Point Label Aware Superpixels [4] on the point label propagation task. This suggests that once the first stage reaches a performance threshold, the second stage performance may saturate.

For the 10 and 25 point label settings, the DeepLabv3+ model trained on masks generated from grid labels exhibits higher pixel accuracy (77.6% and 85.9% for 10 and 25 points respectively) than the HIL generated masks (71.0% and 81.7%). The pixel accuracy calculates the overall correctly classified pixels, which does not take into consideration the class imbalance in the dataset. The mean pixel accuracy and mIoU are calculated as the average per-class scores, and better represent the performance across all classes: the HIL approach outperforms the grid approach by 10.0% and 2.0% for mean pixel accuracy and by 3.4% and 2.1% for mIoU for the 10 point and 25 point cases, respectively (Table II).

TABLE II

PERFORMANCE OF STAGE TWO: SEMANTIC SEGMENTATION WITH DEEPLABV3+ (REFER TO SECTION IV-C FOR METRIC DEFINITIONS), FOR 5 / 10 / 25 / 300 POINT LABELS. ‘F-MSS’ IS FAST MSS [3], ‘PLAS’ IS POINT LABEL AWARE SUPERPIXELS [4], AND ‘D+NN’ IS KNN WITH DENOISED DINOv2 [23] (OURS).

Method	Label Style	PA	mPA	mIoU
		5 / 10 / 25 / 300	5 / 10 / 25 / 300	5 / 10 / 25 / 300
F-MSS	Grid	48.05 / 51.06 / 57.41 / 85.34	5.97 / 10.39 / 14.83 / 63.54	4.41 / 8.56 / 13.04 / 52.16
PLAS - <i>Ens.</i>	Grid	65.73 / 70.18 / 73.60 / 83.72	27.96 / 37.64 / 48.53 / 66.89	19.48 / 26.01 / 36.27 / 52.83
D+NN (Ours)	Grid	72.24 / 77.64 / 85.93 / 85.93	32.16 / 41.61 / 52.59 / 62.99	25.66 / 34.80 / 43.41 / 54.07
D+NN (Ours)	HIL	74.53 / 71.04 / 81.69 / 86.29	41.47 / 51.62 / 54.31 / 63.47	32.96 / 38.21 / 45.46 / 54.62

In these results (Table II), we train and evaluate the models with the “unknown” class included in the dataset. This class was omitted from training and test for the prior approaches [2], [4] because it can contain pixels which belong to one of the “known” classes, thus introducing noise into the training signal. In this case, we aim to establish the efficacy of the different point label propagation approaches by quantifying the performance of models trained on the masks. For the extremely sparse label setting we consider in this work, baseline comparison methods generate propagated ground truth masks dominated by the unknown class (most notably the Fast-MSS approach, as seen in Fig. 6). When these “unknown” pixels are not considered in segmentation task, it leads to model under-fitting which is then misrepresented by the metrics. This effect did not impact results for the original 100-300 point label setting considered in [2], [4], as the larger number of points led to propagated ground truth masks which more closely resembled the ground truth.

B. Ablation Study

1) *Denoising DINOv2 Features*: We employ the denoised version of DINOv2 as described in [23] and highlight its effectiveness as a feature extractor through the results shown in Table III. We present a comparison between the original [22] and denoised [23] DINOv2 deep feature embeddings in Fig. 7. Additionally, we compare the performance of DINOv2 trained with registers [60] to the denoised version trained with registers [23], [60], though this approach did not yield any improvement.

In Fig. 7, we present a visualization of the extracted features with the ground truth mask for each image. We also perform a comparison with a visualization of the CNN features used in the Point Label Aware Superpixel method [4]. Further, we demonstrate that the raw DINOv2

TABLE III
EFFECT OF DIFFERENT DINOv2 FEATURE EXTRACTORS (REFER TO SECTION IV-C FOR METRIC DEFINITIONS)

Denoising	Registers	PA	mPA	mIoU
		5 / 10 / 25 / 300	5 / 10 / 25 / 300	5 / 10 / 25 / 300
✗	✗	68.58 / 73.32 / 76.94 / 88.10	60.23 / 68.04 / 70.97 / 85.58	50.28 / 55.76 / 61.61 / 83.79
✗	✓	68.49 / 73.12 / 76.65 / 87.41	59.79 / 67.48 / 72.44 / 84.84	49.80 / 55.96 / 61.46 / 82.68
✓	✓	70.15 / 75.41 / 78.88 / 88.16	61.85 / 70.75 / 75.81 / 85.42	52.36 / 59.47 / 67.28 / 83.68
✓	✗	71.57 / 76.38 / 80.71 / 89.61	61.46 / 69.87 / 75.91 / 86.45	52.60 / 59.48 / 67.97 / 85.00

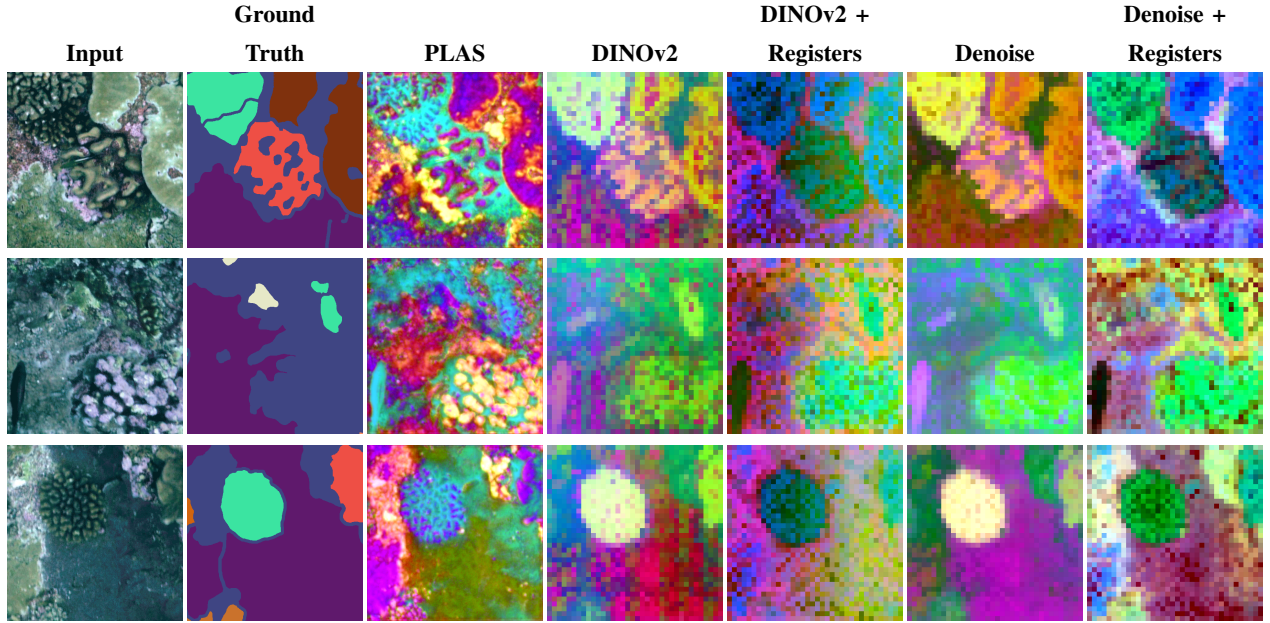


Fig. 7. A comparison is made between Point Label Aware Superpixels (PLAS) [4] features, raw DINOv2 features [22], DINOv2 features trained with registers [60], denoised DINOv2 features [23], and denoised DINOv2 features trained with registers [23], [60] for images in the UCSD Mosaics dataset. In the case of the transformer-based methods, features for each 14x14 pixel patch in the original image are upsampled with bilinear interpolation. Features are visualized in RGB by reducing them with Principal Component Analysis (PCA), where the first three components represent the R, G and B channels. Pixels visualized in similar colors indicate similarity in the deep embedding space. The CNN features used by Point Label Aware Superpixels (PLAS) [4] fail to cluster pixels into distinct segments which align with the expected ground truth segments. The denoised model significantly reduces artifacts from position embeddings, leading to smoother, cleaner features and better clustering performance. In the coral imagery context, the models trained with registers do not seem to improve the feature space.

features contain artifacts due to the way that DINOv2 is trained with position encoding. These artifacts hinder clustering because there can be multiple individuals of the same species spatially separated within the same image. The impact of the positional embeddings was isolated by visualizing the ViT features for a constant value image, both with and without concatenating

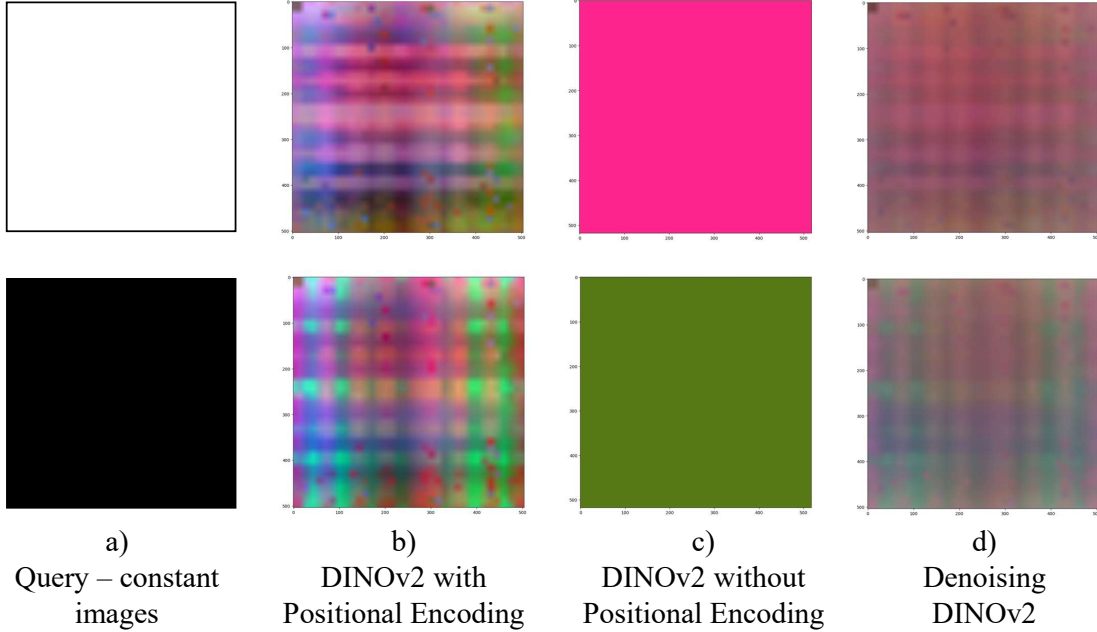


Fig. 8. Features extracted from ViTs are subjected to noisy grid-like artifacts caused by the positional encoding in the architecture. Patch tokens are extracted from the DINOv2 [22] pre-trained ViT, spatially upsampled with bilinear interpolation and then visualized with random projection into the RGB colorspace. Sub-figure a) shows the query images, which are deliberately constant black/white images to demonstrate the artifacts, b) shows the features with noise artifacts from the positional encoding, c) shows that the features are free from noise when the positional encoding is not used, and d) shows that the Denoising ViT [23] significantly reduces the noise artifacts. Positional encoding is needed during training to inform the spatial relationships between the image patches, enabling the model to effectively learn structure and context in images.

the positional embeddings, as seen in Fig. 8. The positional embeddings are necessary during training to encode the relative positions of image patches within the original image [44].

Training DINOv2 with registers reduces some feature artifacts, but not as effectively as the denoising process. The features obtained through training DINOv2 with registers [60] and denoising the features [23] are not as clean as those from the denoised original DINOv2. This is evident in the quantitative results shown in Table III, where the denoised DINOv2 model achieves the highest performance across the three metrics. The denoised version of DINOv2 minimizes the artifacts, yielding a cleaner feature space and therefore improving the point propagation accuracy (Table III).

2) *Weighting the Probability Maps (λ):* We assess how the λ weight, which adjusts the significance of the cosine similarity map (Eq. 2) relative to the distance map (Eq. 3), affects

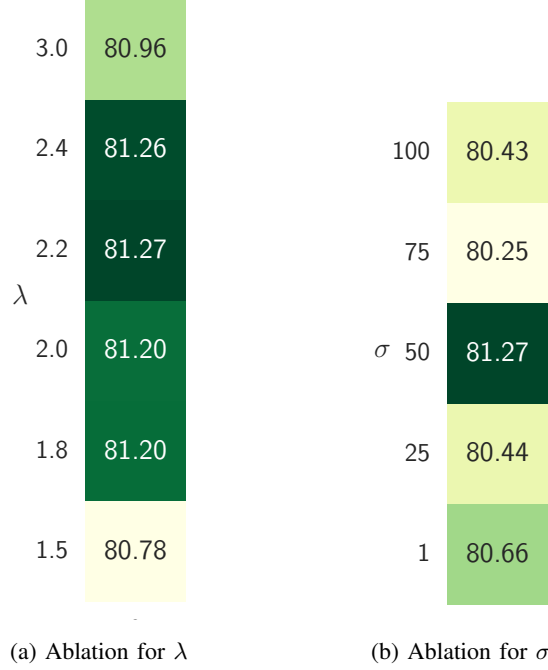


Fig. 9. Label propagation pixel accuracy with 25 point labels, showing that our human-in-the-loop point proposal method is resilient to variation in the value of λ and σ . (a) When the feature similarity map is given more weight, there is a slight improvement in accuracy (we select $\lambda = 2.2$). (b) The highest pixel accuracy is achieved if $\sigma = 50$ for the Gaussian smoothing of the distance map.

performance. As illustrated in Fig. 9, a λ value of 2.2 yields the highest pixel accuracy. However, our method is relatively insensitive to the precise value of λ . The analysis indicates that variations in pixel accuracy are less than 1% across a broad range of λ values tested ($1.5 \leq \lambda \leq 3$).

3) *Exclusion Distance (σ)*: The proposed human-in-the-loop labeling regime accounts for the proximity of labeled pixels by using a Gaussian-smoothed distance mask that measures the distance from all pixels to those that have already been labeled. The Gaussian smoothing introduces the σ hyperparameter, which dictates how near new labeled points can be to existing point labels (Eq. 4). Fig. 9 presents an ablation study on the effect of the σ hyperparameter on the performance of point label propagation, and we find that our method remains effective across various values of σ , ranging from 1 to 100.

4) *Effect of k in KNN*: In this ablation study, we thoroughly evaluate different values of k while varying the number of point labels (5, 10, 25, 100, and 300). The results, displayed in Fig. 10, show that the best performance is achieved with $k = 1$, which corresponds to using a nearest neighbor classifier. This effect is particularly notable with smaller numbers of point labels, as there are fewer examples available for the clustering algorithm. For instance, with only

k	5	41.13	49.34	72.26	82.91	88.57
	3	46.40	57.73	74.94	83.61	88.96
	1	71.56	76.38	81.27	85.60	89.61
		5	10	25	100	300
		Number of Point Labels				

Fig. 10. Ablation for k . Irrespective of the quantity of points labeled per image, the value $k = 1$ always performs the best.

5 point labels in an image, it is likely that there is just one labeled point per class, making it ineffective to consider the majority of three or five neighbors, as only one neighbor will correctly represent the class label.

5) *Effect of Initial Human-Labeled Points:* Our human-in-the-loop labeling regime starts with 10 points labeled centrally within the largest coral instances. Fig. 11 explores how varying the number of initial points labeled by a domain expert affects the pixel accuracy of the augmented ground truth masks. We find that increasing the number of initial labels improves the accuracy of point label propagation. However, if minimizing initial labels is desired, similar results can still be achieved with fewer points: using only 3 initial human-labeled points leads to a decrease in pixel accuracy of 4.5%, 4.8%, and 5.2% when there are 5, 10, and 25 total points, respectively.

C. Effect of Point Label Quantity

Increasing the number of point labels led to better performance in the point label propagation task (Fig. 5). Having access to enough point labels is particularly important for the Fast MSS [3] method. If an image has been annotated with grid-spaced points, Fast MSS shows a significant improvement in mIoU, rising from 7.5% to 86.4% as the number of labels increases from 5 to 300, a difference of 78.9%. In comparison, Point Label Aware Superpixels [4] and our DINOv2 and KNN methods show improvements of 56.5% and 47.1%, respectively (Table I).

While all methods benefit from an increase in point labels, the rate of improvement diminishes as the number of labels grows from 100 to 300. For instance, when increasing grid labels from 100 to 300 points per image, the Fast MSS [3] approach improves by 11.7% in mIoU, compared to a 67.3% improvement when increasing from 5 to 100 points. Similarly, the Point Label Aware

	5	10	25	300
10		76.38	81.27	89.61
5	71.56	74.84	77.67	89.61
3	67.08	71.63	76.08	89.62
1	58.15	66.39	74.71	89.62
	5	10	25	300
	Total Number of Point Labels			

Fig. 11. Ablation for the number of points labeled initially by the domain expert. Increasing the number of points labeled by the expert results in improved performance; here we choose 10 points, although smaller values still result in comparable performance.

Superpixels show a 56.4% improvement in mIoU when moving from 5 to 100 grid points and a 9.5% improvement from 100 to 300 points. In the case of the denoised DINOv2 and KNN, the mIoU improves by 26.2% from 5 to 100 HIL points and by 6.2% from 100 to 300 HIL points.

D. Effect of Point Label Placement Style

All the methods evaluated show advantages when using grid placement of point labels compared to random placement (Fig. 5). This effect is especially notable with the multi-level superpixels (Fast MSS) [3], which shows significant absolute improvements in mIoU with grid-spaced labels over random labels: 13.3%, 9.9%, 6.8% and 6.3% for 50, 100, 200 and 300 points, respectively. The Point Label Superpixels also benefit from grid spacing, with improvements of 8.4%, 4.8%, 4.4% and 3.9% for the same label quantities. Grid-spaced labels ensure uniform coverage across the entire image and make optimal use of each label. As illustrated in Fig. 6, randomly placed labels can be clustered closely together, diminishing the amount of useful information.

Fig. 6 also shows that, with very few point labels (5 to 10 per image), there is considerable benefit from utilizing domain expert knowledge to select points centrally within instances. Further pixels can then be iteratively selected using the point propagation model described in Section III. The augmented ground truth masks produced by our approach (top two rows of Fig. 6) are

significantly closer to the actual ground truth compared to previous methods. We plan to explore applying our human-in-the-loop labeling regime to other techniques in future research.

When labels are sparse, multi-level superpixel methods [2], [3] struggle because they depend on layering labeled regions from various scales. Similarly, the point label superpixel method faces challenges in sparse cases, as its superpixel boundaries are not forced to align with instance boundaries by the conflicting point labels [4]. Our method performs well with sparse labels because it assigns the correct class to pixels even if they are spatially distant from labeled points, by relying on clustering in the deep feature space.

VI. CONCLUSION

This study has shown that the general foundation model, DINOv2, can be effectively utilized for point label propagation in underwater imagery without requiring fine-tuning. By leveraging denoised DINOv2 features and a straightforward KNN algorithm, we generate augmented ground truth masks. When combined with our human-in-the-loop labeling approach, we achieve significant improvements in mIoU: 19.7%, 18.3%, and 8.2% for 5, 10, and 25 point labels, respectively, when compared to previous state-of-the-art methods. Even when using DINOv2 features and KNN with grid-spaced point labels, we outperform earlier methods on the UCSD dataset with improvements of 5.8%, 10.0%, and 4.6% mIoU for 5, 10, and 25 point labels, respectively.

These performance gains are evident in the semantic segmentation task as well. Training a DeepLabv3+ model on augmented ground truth masks created with DINOv2, KNN, and our human-in-the-loop labeling approach yields improvements of 8.8% in pixel accuracy and 13.5% in mIoU with just 5 point labels. If the human-in-the-loop regime is not employed, the semantic segmentation still benefits, showing an 6.5% increase in pixel accuracy and a 6.2% increase in mIoU.

Our comprehensive studies on the impact of the number of point labels and their placement reveal that grid labels consistently enhance point label propagation accuracy compared to random labels. For more than 100 points per image, improvements in point label propagation are minimal with either the Point Aware Superpixels [4] method or our DINOv2 and KNN approach. However, for very few points per image, domain expert input in selecting which pixels to label proves advantageous. This work highlights the effectiveness of general foundation models for complex,

domain-specific tasks and significantly boosts performance and annotation efficiency in scenarios with extremely sparse labels.

Our method could be extended to treat the human-in-the-loop scenario as an active learning problem during deployment. While our current human-in-the-loop framework is focused on point label propagation, it could be modified to detect novel species in real-time and utilize domain expert input to incorporate new classes through active learning. In this scenario, a deep learning model could be deployed on a robotic platform and used to predict previously unseen targets in the footage. The human-in-the-loop scheme would then present identified anomalies to a domain expert in real-time. The expert determines whether the class is of interest and if so, specifies the class label. As further examples are detected, the examples could be incorporated in an active learning framework where the model iteratively incorporates additional classes. Underwater anomaly detection could have significant implications for species discovery and would enable models to be quickly adapted to new locations and species encountered, rather than exhibiting degraded performance when deployed in a different setting from the training data.

Another future direction could be to obtain features usable across different images. The architecture and training of the DINOv2 foundation model incorporates global image information into per-pixel deep features, resulting in deep features of the same coral species from different images not being similar in the deep embedding space. Enabling feature similarity across different images could eliminate the need for the second stage of this architecture, thereby simplifying the framework by removing the need for DeepLabv3+ training.

The human-in-the-loop labeling regime was simulated using the ground truth masks with the UCSD Mosaics dataset [2], [51]. Future research might include testing with multiple domain experts to assess how experts interact with the human-in-the-loop regime and how the inter-observer variability affects performance.

Given the increasing frequency of global coral bleaching events, there is a growing need for detailed monitoring of coral reefs and the identification of temperature-resistant species. Future research could expand our approach to enable precise recognition of coral reef health indicators.

The large number of coral species, the complexity of coral imagery, and the location-specific nature of coral visual features mean that real-world deployment of multi-species image analysis systems must efficiently gather data in the target environment, label a few points, and then train and deploy a model. With an estimated 800 species of hard corals alone [61], collecting

and annotating sufficient data for a global model is impractical. Instead, we believe that efforts should focus on rapidly training and deploying models tailored to specific locations. The method proposed in this work contributes to this goal by significantly reducing the amount of required annotations while achieving accurate point label propagation.

This work has focused on segmentation of coral imagery, however the approach presented could be used for broad-scale surveys of other species of interest, including seagrass meadows and algae. Further, the Stony Coral Tissue Loss Disease is an aggressive disease which can kill entire colonies of stony corals in a period of months [62]. The spread of the disease is not fully understood so underwater monitoring over varying spatial and temporal scales is important for detection and management of the disease [62], [63]. As there is limited data available which captures the visual characteristics of the disease, the human-in-the-loop labeling regime presented in this work could be used to enable fast annotation, training and deployment of disease detection models.

In addition to marine surveys, environmental monitoring on land is often subject to similar data availability and annotation challenges. Broad-scale surveys completed using Unmanned Aerial Vehicles (UAVs) or satellite imagery require annotation for training deep learning models, and in domain-specific applications this can be expensive and time-consuming, as it is for underwater imagery. Future work could investigate the application of our presented human-in-the-loop labeling approach on aerial and remote data for domain-specific environmental monitoring tasks.

ACKNOWLEDGMENTS

This work was done in collaboration between QUT and CSIRO Data61. S.R., F.M., N.S., and T.F. acknowledge continued support from the Queensland University of Technology (QUT) through the Centre for Robotics. T.F. acknowledges funding from Intel Labs via grant RV3.290.Fischer and an ARC Laureate Fellowship FL210100156 to Prof. Michael Milford. Computational resources and services used in this work were provided by the eResearch Office, Queensland University of Technology, Brisbane, Australia.

REFERENCES

- [1] I. Alonso and A. C. Murillo, "Semantic segmentation from sparse labeling using multi-level superpixels," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2018.

- [2] I. Alonso, M. Yuval, G. Eyal, T. Treibitz, and A. C. Murillo, "CoralSeg: Learning coral segmentation from sparse annotations," *Journal of Field Robotics*, vol. 36, 2019.
- [3] J. P. Pierce, Y. Rzhannov, K. Lowell, and J. A. Dijkstra, "Reducing annotation times: Semantic segmentation of coral reef survey images," in *Proceedings of the OCEANS Conference*, 2020.
- [4] S. Raine, R. Marchant, B. Kusy, F. Maire, and T. Fischer, "Point label aware superpixels for multi-species segmentation of underwater imagery," *IEEE Robotics and Automation Letters*, vol. 7, 2022.
- [5] E. M. Ditría, C. A. Buelow, M. Gonzalez-Rivero, and R. M. Connolly, "Artificial intelligence and automated monitoring for assisting conservation of marine ecosystems: A perspective," *Frontiers in Marine Science*, vol. 9, 2022.
- [6] M. Dunbabin, J. Manley, and P. L. Harrison, "Uncrewed maritime systems for coral reef conservation," in *Proceedings of the OCEANS Conference*, 2020.
- [7] D. Gregorek, A. Tibebe, E. Caudet, C. Barrera, and R. Bachmayer, "Long-endurance optical seafloor imaging using underwater gliders: Concept, development and initial trials," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2023.
- [8] K. Mizuno, K. Terayama, S. Tabeta, S. Sakamoto, Y. Matsumoto, Y. Sugimoto, T. Ogawa, K. Sugimoto, H. Fukami, M. Sakagami *et al.*, "Development of an efficient coral-coverage estimation method using a towed optical camera array system [Speedy Sea Scanner (SSS)] and deep-learning-based segmentation: A sea trial at the Kujuku-Shima Islands," *IEEE Journal of Oceanic Engineering*, vol. 45, 2019.
- [9] A. Mahmood, M. Bennamoun, S. An, F. A. Sohel, F. Boussaid, R. Hovey, G. A. Kendrick, and R. B. Fisher, "Deep image representations for coral image classification," *IEEE Journal of Oceanic Engineering*, vol. 44, 2018.
- [10] L. Xu, M. Bennamoun, S. An, F. Sohel, and F. Boussaid, "Deep learning for marine species recognition," in *Handbook of Deep Learning Applications*, 2019.
- [11] H. Runyan, V. Petrovic, C. B. Edwards, N. Pedersen, E. Alcantar, F. Kuester, and S. A. Sandin, "Automated 2D, 2.5D, and 3D segmentation of coral reef pointclouds and orthoprojections," *Frontiers in Robotics and AI*, vol. 9, 2022.
- [12] Y. Guo, Y. Liu, T. Georgiou, and M. S. Lew, "A review of semantic segmentation using deep neural networks," *International Journal of Multimedia Information Retrieval*, vol. 7, 2018.
- [13] M. Zhang, Y. Zhou, J. Zhao, Y. Man, B. Liu, and R. Yao, "A survey of semi-and weakly supervised semantic segmentation of images," *Artificial Intelligence Review*, 2019.
- [14] Z. Fu, R. Chen, Y. Huang, E. Cheng, X. Ding, and K.-K. Ma, "MASNet: A robust deep marine animal segmentation network," *IEEE Journal of Oceanic Engineering*, 2023.
- [15] M. Li, H. Zhang, A. Gruen, and D. Li, "A survey on underwater coral image segmentation based on deep learning," *Geo-spatial Information Science*, 2024.
- [16] X. Sun, J. Shi, L. Liu, J. Dong, C. Plant, X. Wang, and H. Zhou, "Transferring deep knowledge for object recognition in low-quality underwater videos," *Neurocomputing*, vol. 275, 2018.
- [17] L. Jin and H. Liang, "Deep learning for underwater image recognition in small sample size situations," in *Proceedings of the OCEANS Conference*, 2017.
- [18] K. E. Kohler and S. M. Gill, "Coral Point Count with Excel extensions (CPCe): A Visual Basic program for the determination of coral and substrate coverage using random point count methodology," *Computers & Geosciences*, vol. 32, 2006.
- [19] M. González-Rivero *et al.*, "Monitoring of coral reefs using artificial intelligence: A feasible and cost-effective approach," *Remote Sensing*, vol. 12, 2020.
- [20] O. Beijbom, P. J. Edmunds, D. I. Kline, B. G. Mitchell, and D. Kriegman, "Automated annotation of coral reef survey images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2012.

- [21] M. González-Rivero, P. Bongaerts, O. Beijbom, O. Pizarro, A. Friedman, A. Rodriguez-Ramirez, B. Upcroft, D. Laffoley, D. Kline, C. Bailhache *et al.*, “The Catlin seaview survey—kilometre-scale seascape assessment, and monitoring of coral reef ecosystems,” *Aquatic Conservation: Marine and Freshwater Ecosystems*, vol. 24, 2014.
- [22] M. Oquab, T. Darcet, T. Moutakanni, H. V. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang, S.-W. Li, I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski, “DINOv2: Learning robust visual features without supervision,” *Transactions on Machine Learning Research*, 2024.
- [23] J. Yang, K. Z. Luo, J. Li, K. Q. Weinberger, Y. Tian, and Y. Wang, “Denoising vision transformers,” *Proceedings of the European Conference on Computer Vision*, 2024.
- [24] Y. Li, J. Liu, B. Kusy, R. Marchant, B. Do, T. Merz, J. Crosswell, A. Steven, L. Tychsen-Smith, D. Ahmedt-Aristizabal *et al.*, “A real-time edge-AI system for reef surveys,” in *Proceedings of the Annual International Conference on Mobile Computing And Networking*, 2022.
- [25] S. Mou, D. Tsai, and M. Dunbabin, “Reconfigurable robots for scaling reef restoration,” *arXiv preprint arXiv:2205.04612*, 2022.
- [26] M. González-Rivero, O. Beijbom, A. Rodriguez-Ramirez, T. Holtrop, Y. González-Marrero, A. Ganase, C. Roelfsema, S. Phinn, and O. Hoegh-Guldberg, “Scaling up ecological measurements of coral reefs using semi-automated field image collection and analysis,” *Remote Sensing*, vol. 8, 2016.
- [27] Z. Ziqiang, X. Yaofeng, L. Haixin, Y. Zhibin, and S.-K. Yeung, “CoralVOS: Dataset and benchmark for coral video segmentation,” *arXiv preprint arXiv:2310.01946*, 2023.
- [28] H. Zhang, M. Li, J. Zhong, and J. Qin, “CNet: A novel seabed coral reef image segmentation approach based on deep learning,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024.
- [29] H. Zhang, A. Grün, and M. Li, “Deep learning for semantic segmentation of coral images in underwater photogrammetry,” *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 2, 2022.
- [30] Y. W. Sui, K. X. Ming, M. Meghjani, N. Raghavan, C. Jegourel, and K. Kang, “An automated data processing pipeline for coral reef monitoring,” in *Proceedings of the OCEANS Conference*, 2022.
- [31] J. Zhong, M. Li, H. Zhang, and J. Qin, “Combining photogrammetric computer vision and semantic segmentation for fine-grained understanding of coral reef growth under climate change,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023.
- [32] D. P. Furtado, E. A. Vieira, W. F. Nascimento, K. Y. Inagaki, J. Bleuel, M. A. Z. Alves, G. O. Longo, and L. S. Oliveira, “#DeOlhoNosCorais: A polygonal annotated dataset to optimize coral monitoring,” *PeerJ*, vol. 11, 2023.
- [33] M. J. Islam, C. Edge, Y. Xiao, P. Luo, M. Mehtaz, C. Morse, S. S. Enan, and J. Sattar, “Semantic segmentation of underwater imagery: Dataset and benchmark,” in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2020.
- [34] D. Langenkämper, M. Zurowietz, T. Schoening, and T. W. Nattkemper, “Biigle 2.0-browsing and annotating large marine image collections,” *Frontiers in Marine Science*, vol. 4, 2017.
- [35] O. Beijbom, P. J. Edmunds, C. Roelfsema, J. Smith, D. I. Kline, B. P. Neal, M. J. Dunlap, V. Moriarty, T.-Y. Fan, C.-J. Tan *et al.*, “Towards automated annotation of benthic survey images: Variability of human experts and operational modes of automation,” *PloS One*, vol. 10, 2015.
- [36] M. Zurowietz, D. Langenkämper, B. Hosking, H. A. Ruhl, and T. W. Nattkemper, “MAIA: A machine learning assisted image annotation method for environmental monitoring and exploration,” *PloS One*, vol. 13, 2018.
- [37] I. Urbina-Barreto, R. Garnier, S. Elise, R. Pinel, P. Dumas, V. Mahamadaly, M. Facon, S. Bureau, C. Peignon, J.-P. Quod

- et al.*, “Which method for which purpose? A comparison of line intercept transect and underwater photogrammetry methods for coral reef surveys,” *Frontiers in Marine Science*, vol. 8, 2021.
- [38] G. Pavoni, M. Corsini, F. Ponchio, A. Muntoni, C. Edwards, N. Pedersen, S. Sandin, and P. Cignoni, “TagLab: AI-assisted annotation for the fast and accurate semantic segmentation of coral reef orthoimages,” *Journal of Field Robotics*, vol. 39, 2022.
 - [39] Z. Zhang, P. Kaveti, H. Singh, A. Powell, E. Fruh, and M. E. Clarke, “An iterative labeling method for annotating marine life imagery,” *Frontiers in Marine Science*, vol. 10, 2023.
 - [40] X. Yu, B. Ouyang, J. C. Principe, S. Farrington, J. Reed, and Y. Li, “Weakly supervised learning of point-level annotation for coral image segmentation,” in *Proceedings of the OCEANS Conference*, 2019.
 - [41] X. Yu, Y. Ma, S. Farrington, J. Reed, B. Ouyang, and J. C. Principe, “Fast segmentation for large and sparsely labeled coral images,” in *Proceedings of the IEEE International Joint Conference on Neural Networks*, 2019.
 - [42] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, “Segment anything,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
 - [43] X. Zou, J. Yang, H. Zhang, F. Li, L. Li, J. Wang, L. Wang, J. Gao, and Y. J. Lee, “Segment everything everywhere all at once,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
 - [44] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *Proceedings of the International Conference on Learning Representations*, 2021.
 - [45] R. Zhang, Z. Jiang, Z. Guo, S. Yan, J. Pan, H. Dong, P. Gao, and H. Li, “Personalize segment anything model with one shot,” in *Proceedings of the International Conference on Learning Representations*, 2024.
 - [46] F. Chen, M. V. Giuffrida, and S. A. Tsaftaris, “Adapting vision foundation models for plant phenotyping,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
 - [47] M. Baharoon, W. Qureshi, J. Ouyang, Y. Xu, K. Phol, A. Aljouie, and W. Peng, “Towards general purpose vision foundation models for medical image analysis: An experimental study of DINOv2 on radiology benchmarks,” *arXiv preprint arXiv:2312.02366*, 2023.
 - [48] O. Siméoni, G. Puy, H. V. Vo, S. Roburin, S. Gidaris, A. Bursuc, P. Pérez, R. Marlet, and J. Ponce, “Localizing objects with self-supervised transformers and no labels,” in *Proceedings of the British Machine Vision Conference*, 2021.
 - [49] J. P. Huix, A. R. Ganeshan, J. F. Haslum, M. Söderberg, C. Matsoukas, and K. Smith, “Are natural domain foundation models useful for medical image classification?” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024.
 - [50] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
 - [51] C. B. Edwards *et al.*, “Large-area imaging reveals biologically driven non-random spatial patterns of corals at a remote reef,” *Coral Reefs*, vol. 36, 2017.
 - [52] E. Mosqueira-Rey, E. Hernández-Pereira, D. Alonso-Ríos, J. Bobes-Bascarán, and Á. Fernández-Leal, “Human-in-the-loop machine learning: A state of the art,” *Artificial Intelligence Review*, vol. 56, 2023.
 - [53] L. Jiang, S. Liu, and C. Chen, “Recent research advances on interactive machine learning,” *Journal of Visualization*, vol. 22, 2019.
 - [54] Q. Chen, O. Beijbom, S. Chan, J. Bouwmeester, and D. Kriegman, “A new deep learning engine for CoralNet,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.

- [55] B. Kellenberger, D. Tuia, and D. Morris, “AIDE: Accelerating image-based ecological surveys with interactive machine learning,” *Methods in Ecology and Evolution*, vol. 11, 2020.
- [56] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, “Encoder-decoder with atrous separable convolution for semantic image segmentation,” in *Proceedings of the European Conference on Computer Vision*, 2018.
- [57] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *Proceedings of the Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [58] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with GPUs,” *IEEE Transactions on Big Data*, vol. 7, 2019.
- [59] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, P. Martinez-Gonzalez, and J. Garcia-Rodriguez, “A survey on deep learning techniques for image and video semantic segmentation,” *Applied Soft Computing*, vol. 70, 2018.
- [60] T. Darcet, M. Oquab, J. Mairal, and P. Bojanowski, “Vision transformers need registers,” in *Proceedings of the International Conference on Learning Representations*, 2024.
- [61] A. Dietzel, M. Bode, S. R. Connolly, and T. P. Hughes, “The population sizes and global extinction risk of reef-building coral species at biogeographic scales,” *Nature Ecology and Evolution*, vol. 5, 2021.
- [62] E. Papke, A. Carreiro, C. Dennison, J. M. Deutsch, L. M. Isma, S. S. Meiling, A. M. Rossin, A. C. Baker, M. E. Brandt, N. Garg *et al.*, “Stony coral tissue loss disease: a review of emergence, impacts, etiology, diagnostics, and intervention,” *Frontiers in Marine Science*, vol. 10, 2024.
- [63] I. R. Combs, M. S. Studivan, R. J. Eckert, and J. D. Voss, “Quantifying impacts of stony coral tissue loss disease on corals in southeast florida through surveys and 3D photogrammetry,” *PLOS One*, vol. 16, 2021.