



**HAL**  
open science

# Towards a Fully Automated Underwater Census for Fish Assemblages in the Mediterranean Sea

Kilian Bürgi, Charles Bouveyron, Diane Lingrand, Benoit Dérijard, Frédéric Precioso, Cécile Sabourault

## ► To cite this version:

Kilian Bürgi, Charles Bouveyron, Diane Lingrand, Benoit Dérijard, Frédéric Precioso, et al.. Towards a Fully Automated Underwater Census for Fish Assemblages in the Mediterranean Sea. 2024. hal-04690514

**HAL Id: hal-04690514**

**<https://hal.science/hal-04690514>**

Preprint submitted on 6 Sep 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Towards a Fully Automated Underwater Census for Fish Assemblages in the Mediterranean Sea

Kilian Bürgi<sup>a,b</sup>, Charles Bouveyron<sup>b</sup>, Diane Lingrand<sup>c</sup>, Benoit Dérijard<sup>a</sup>,  
Frédéric Precioso<sup>c</sup>, and Cécile Sabourault<sup>a,\*</sup>

\*Corresponding author, Cecile.SABOURAULT@univ-cotedazur.fr

<sup>a</sup>Université Côte d’Azur, CNRS, ECOSEAS, Nice, France

<sup>b</sup>Université Côte d’Azur, Inria, CNRS, Laboratoire J.A.Dieudonné, Maasai team, Nice, France

<sup>c</sup>Université Côte d’Azur, Inria, CNRS, I3S, Maasai team, Nice, France

*{Kilian.BURGI,Charles.BOUVEYRON,Diane.LINGRAND}@univ-cotedazur.fr*  
*{Benoit.DERIJARD,Frederic.PRECIOSO,Cecile.SABOURAULT}@univ-cotedazur.fr*

June 6, 2024

## Abstract

Assessing underwater biodiversity is a labour-intensive and costly procedure whilst being crucial to measure the extent of local fish stock declines. In most cases, Underwater Visual Census (UVC) is the method of preference, however this can be human-costly and is limited by meteorological and logistic factors. Advances in technology allows the utilisation of more autonomous video recording methods (*i.e.* Remote Operated Vehicles (ROV)) which work around the aforementioned limitations. This study used a transect-wise UVC coupled with diver operated videos (DOV) simulating an ROV. For the video analysis, a comprehensive fully automated pipeline was developed to extract frames from DOV and perform color correction. This pipeline integrates a YOLO-based model for the detection of 20 Mediterranean fish species validating presence or absence of each species within individual transect. This study was conducted to evaluate the feasibility of utilising video-based methods for UVC with minimal human-dependence. The automation of the video analysis showed accordance with the manual video counting enabling an autonomous and bias-free procedure for video assessment. In conclusion, utilising a minimal-human-dependent video method disconnects the data acquisition from limiting factor (*i.e.* meteorological and logistic) and automation of this video analysis will significantly reduce the labour and time required by researchers. For future

fieldwork campaigns, the video data collection protocol needs to be adjusted to better resemble the traditional UVC and bring forward this acquisition method.

Diver operated video, Automated UVC, Deep learning, Object detection, Marine biology, Marine protected areas

## Highlights

1. Applying YOLOv7 deep learning model on diver-operated video (DOV) transects & sites.
2. UVC & video methods (manual & automated) were compared identifying presence/absence.
3. UVC & video data combination shows more complete site evaluation.
4. Total of 85% of the species *in situ* were correctly identified by automation.
5. *Epinephelus marginatus* was seen twice as much inside of MPA than outside.

## 1 Introduction

The marine environment is facing multiple stressors that have a significant impact on their ecosystems (Gissi *et al.*, 2021). Artificialization of shorelines (Carranza *et al.*, 2019), overfishing (Demirel *et al.*, 2020), masstourism (Mejjad *et al.*, 2022) and climate change (Doney *et al.*, 2012; Smale *et al.*, 2019) are impacting the fish communities all over the world and especially in high touristic areas such as the French Riviera. High demand of fish meat as protein source led and still leads to fish populations declining up to 99% for very extreme sites in the past decades (Myers *et al.*, 1997; Vasilakopoulos *et al.*, 2014). There is a developmental need to survey, monitor and conserve these fragile marine areas and gather information on the current state in which they are in. For this purpose marine protected areas (MPA) are established and function as safe havens for fish populations to recover and proliferate. Inside of an MPA, human activities such as fishing, diving or anchoring are limited or prohibited to ensure the decompression of the ecosystem. However the efficient management of these MPA is complicated and underlies a careful and holistic procedure of data collection to make decisions regarding its policies.

This data collection involves both, abiotic factors (*i.e.* bathymetric and physical data) and biotic factors such as biodiversity. Whilst abiotic factors can be measured in real-time with corresponding probes, the measure of biotic factors are either invasive (*i.e.* experimental fishing, catch and release or other methods involving the landing of the fish)

or less-invasive but very cost- and labour-intensive (*i.e.* underwater visual census (UVC) or video-based methods). The second approach of lesser-invasion should be considered first since the health of the ecosystem should not be compromised to evaluate it.

UVC has been standardised for decades (Kulbicki and Sarramagna, 1999) and has the advantages to be well established. It forms the best global coverage of biodiversity assessment (Caldwell *et al.*, 2016) but it suffers from diver bias that can - in high biodiversity regions - lead up to a 25% over- or under-estimation of abundance and species richness (Mcclanahan *et al.*, 2007). Alongside this bias, UVC is high-cost, requires a lot of expert knowledge and is heavily dependent on weather and logistics due to safety reasons for the divers. To counteract these factors unmanned recording techniques have risen. Most of the existing video assisted studies use baited remote underwater videos (BRUV) which are temporary stationary cameras with bait attached in front of it to lure fish in front of the camera to evaluate the biodiversity or other ecologically relevant information (Mclean *et al.*, 2005). The methodology of BRUV is less-invasive and give a good insight in terms of biodiversity but are susceptible to wrong decisions in choosing the right bait and location which leads to a bias towards bait-preferential fish.

To gain independence of mentioned factors, a traditional UVC was coupled with video assistance (Grorud-Colvert *et al.*, 2021) to help the diver return visually to the transect in case of need of confirmation and evaluation of complexity or coverage of the present transect substrate. The video recordings following the divers point of view are hypothesized to simulate a remote operated vehicle (ROV) and capture a similar biology as the diver itself. With this technique, numerous videos were recorded and looking through them was time intensive. In the era of digitalisation, artificial intelligence and its subdisciplines are a vastly rising field of research in marine ecology (Malde *et al.*, 2020; Rubbens *et al.*, 2023) and could give greater insights on the images and videos recorded in many different contexts (Vabø *et al.*, 2021; Wu *et al.*, 2022). Deep learning (DL) has shown promising results in the analysis of underwater imagery (*i.e.* Spampinato *et al.*, 2016; Jalal *et al.*, 2020). For example achieving excellent results in in-trawl images with mean average precision (mAP) of 0.845 (Allken *et al.*, 2021a) with a creatively acquired and expanded dataset (Allken *et al.*, 2021b) is truly inspiring. A study closer related to the proposed study is for example Xu and Matzner, 2018. This team achieved an mAP of 0.54 on a real life applicable dataset while Knausgård *et al.*, 2022 achieved an mAP of 0.84 on a temperate fish dataset which helped a lot to benchmark and elevate this study. However, these studies evaluate videos derived from stationary cameras and focus on machine learning metrics (F1 score, recall, precision and mAP), leaving out the biological aspect of it. Whilst the importance of the machine learning performance metrics is undeniable, in real life applications the biological or methodological metrics need to be incorporated to create the bigger picture

of the study and not only focus on one aspect of it. Connolly *et al.*, 2022 used ROV video data and is automatically detecting two economically important species on the Australian coast to evaluate moving videos in comparison to stationary cameras. This study achieved very good results to count fish in frames but showed differences in model performance for the two different fish species.

In the proposed study the focus was on a video-assisted UVC protocol and how DL can help hastening the diver observations and assist the presence/absence video transect analysis. The framework of You Only Look Once (YOLO) is widely used in marine computer vision (Mohamed *et al.*, 2020; Park and Kang, 2020; Priyankan and Fernando, 2021; Muksit *et al.*, 2022) and the version 7 (C.-Y. Wang *et al.*, 2023) was used in this study as the DL model for fish detection. The model was evaluated in a first experiment in its performance capabilities in detecting 20 classes (19 most prominent local fish species and 1 'Other' class). In a second experiment on a different dataset, the same 20 species were detected, transect-wise concatenated and formatted into a presence/absence table for each of the species. Manual video counts performed by a marine biology expert and professional diver-gathered UVC data were compared to the presence/absence table derived from the detections. This study seeks to compare how manual video data, data generated by artificial intelligence, and data collected by scuba divers differ in their perceptions of the diversity of 20 different marine fish species in the Mediterranean Sea. The aim is to compare these methods and define the degree of differences allowing to make adaption propositions for future fieldwork campaigns.

## 2 Materials and methods

### 2.1 Study area & data collection

The training dataset ( $DATA_T$ ) was gathered in eight different locations of the French Riviera in the Mediterranean Sea and followed the same UVC protocol on each site (Harmelin-Vivien *et al.*, 1985). The depth ranged from 1-37m and was executed during the whole year in 2022 (cold- and warm season) to have the full range of conditions and possibilities of fish occurrences.

The experimental dataset ( $DATA_E$ ) is evaluated in terms of methodology and was recorded in October 2023 in two distinct areas, one no-take zone and one Natura2000 site, which both have elevated biodiversity - 'Cap Roux' and 'Corniche Varoise' in the French Riviera. The specific coordinates & meta data can be found in the supplementary material (Table S1). A total of 64 videos, each corresponding to a transect, from 14 sites (8 on seagrass meadows & 6 on rocky substrates) were evaluated and compared. Each site

consists of 3 to 6 transects depending on the availability of the video recordings and the UVC data from the divers.

For the recording of the videos, GoPro version 9 cameras were used. They were mounted on the clipboard (Fig. 1) where the divers note the number of fish per species with their respective size category (variable category number per species). These videos were recorded with a framerate of 24 frames per second (FPS) and a full high definition resolution (1920x1080px). Frames were extracted from these recordings with a FPS of 1 for  $DATA_T$  and with a FPS of 5 for  $DATA_E$ . A fish less than 1 second (less than 5 frames) in the videos of  $DATA_E$  will not be considered in the methodology evaluation due to the unlikelihood of it being an actual detection.



Figure 1: Example image of the GoPro 9 montage of a diver on a transect that records the current divers point of view. The background shows the measuring tape can be seen to stay on track with the UVC protocol.

## 2.2 Preprocessing of the images & training

To ensure a good species coverage, 19 different species and an 'Other' class were labelled manually in the frames resulting in 13'033 images (131 videos in the training set and 47 independent videos in the test set) in  $DATA_T$  having a total of 68'573 (train = 40'379, test = 28'194) individual fish labels (species breakdown in Table S3) and 8'739 miscellaneous labels such as background and diver. The 'Other' class includes species (Table S2) that have insufficient occurrences in the test videos ( $n < 100$ ).

Since there is wide range of conditions in the videos, a preprocessing was applied to enhance each image colour range. For this purpose a pretrained UIEC<sup>2</sup>-Net model (Y. Wang *et al.*, 2021) was utilised to enhance the images. For the training the YOLOv7 algorithm was used (C.-Y. Wang *et al.*, 2023) and the model was pre-trained on two public fish datasets - DeepFish (Garcia-D'Urso *et al.*, 2022) and OzFish (AIMS *et al.*, 2019) -

and then fine tuned on the data described above. As a preprocessing, images have been rescaled to 960x960px (from the default 640x640px suggested by YOLO). The training was run on a high performance GPU cluster for 150 epochs to guarantee convergence. The best performing weights (maximised mAP@0.5:0.95 on the validation set) were automatically provided by YOLO and chosen as the weights for the DL model.

The whole process from receiving the recordings to the different metric evaluation is depicted in Figure 2. The pipeline depicts the DL detection for the diver after returning from fieldwork.

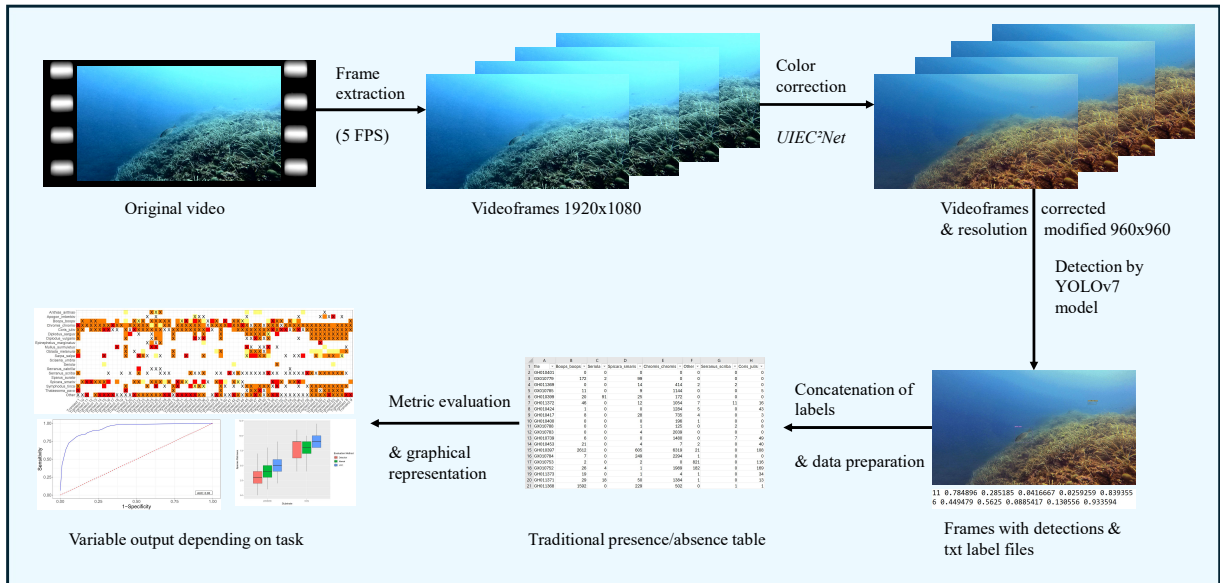


Figure 2: The pipeline of the videos gathered for the method evaluation.

## 2.3 Metrics of Evaluation

### 2.3.1 Deep learning model evaluation

To evaluate the first experiment of the detection capability of the YOLO model, the mean average precision (mAP - Eq. 1) which consists of the area under the curve (AUC) of the precision-recall curve was used. Precision (Eq. 2) indicates the proportion of correct detections among all detections made by the model, while recall (Eq. 3) represents the proportion of actual correct detections that the model successfully identified. These two metrics are common metrics used in DL to evaluate the models performance on a given task. The abbreviations FN, FP, TN and TP stand for False Negatives, False Positives, True Negatives and True Positives respectively.

$$\text{meanAveragePrecision} = mAP = \frac{1}{n} \sum_{k=1}^{k=n} AP_k \quad (1)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

### 2.3.2 Binary evaluation of automated & manual video data

For evaluation of the second experiment with  $DATA_E$ , the transect-wise presence/absence study, the following metrics are calculated for the manual video counting compared to the presence/absence table generated from the detections. The Accuracy (Eq. 4), Sensitivity (Eq. 5), Specificity (Eq. 6), 1-Specificity (Eq. 7), Cohens Kappa (Eq. 8, Landis and Koch, 1977) and true skill statistic (TSS, Eq. 9) are used to get a holistic insight on how the model is performing on this specific task. Cohens Kappa is used to evaluate the inter-reliability between two or more measuring methods - in this case between the manual video count and the DL predictions. The TSS disconnects the Cohens Kappa from its prevalence problem (Allouche *et al.*, 2006).

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad (4)$$

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (5)$$

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (6)$$

$$1 - \text{Specificity} = 1 - \frac{TN}{TN + FP} = \frac{FP}{TN + FP} \quad (7)$$

$$\text{CohensKappa} = \frac{(TP + FP)(TP + FN) + (FN + TN)(TN + FP)}{(TP + FP + TN + FN)^2} \quad (8)$$

$$\text{TSS} = \text{Sensitivity} + \text{Specificity} - 1 \quad (9)$$



### 3 Results

#### 3.1 Deep learning model

For the first experiment, the confusion matrix for the 20 classes is presented. Background FP refers to incorrect detections of the DL model that identify parts of the background as objects, while background FN refer to missed detections against the background. The darker coloration indicates a higher relative classification in this specific true class. This can suggest either higher overall accuracy (when it is on the diagonal, indicating more correct predictions) or higher misclassification rates (when it is off the diagonal, indicating more instances wrongly classified). The confusion matrix (Fig. 3) shows a clear diagonal line for most of the species covered except for the classes *Mullus surmuletus* and *Other*. Overall a mean average precision (mAP) of 0.56 is achieved. The model has an overall precision of 0.66 with 18'590 TP and 9'716 FP. Individual species precision (Table S4) values range from 0.13 (class 'Other') to 0.91 (class 'Sarpa\_salpa'). Second worst after the 'Other' class is 'Seriola' with a precision of 0.41. Recall values (overall 0.65) range from 0.23 for *M. surmuletus* to 0.75 for *Coris julis*.

	background FN	2642	477	121	538	94	49	84	267	51	105	78	222	442	241	154	277	25	233	97	115	0
	Oblada_melanura	50	32	0	23	8	1	0	4	0	6	2	12	1	6	14	36	2	1	0	197	76
	Thalassoma_pavo	0	0	5	1	1	0	9	0	3	0	0	5	3	12	0	1	0	0	101	0	12
	Boops_boops	57	0	2	20	0	0	0	14	5	0	0	2	7	1	255	1	1	1147	2	3	379
	Seriola	0	2	0	108	0	2	0	0	0	0	11	6	0	1	0	6	119	0	0	0	34
	Diplodus_sargus	7	422	0	98	1	3	0	17	2	5	27	29	2	10	2	915	3	0	1	99	199
	Spicara_smaris	27	4	0	1	0	0	2	12	2	0	0	2	0	0	669	0	0	144	0	3	518
	Symphodus_tinca	14	2	24	59	0	0	2	15	13	0	0	30	41	440	0	13	0	1	10	17	143
	Coris_julis	24	3	55	14	0	1	2	16	39	0	0	24	1914	67	5	6	0	1	23	2	439
	Other	49	18	31	377	0	2	5	5	4	19	2	145	16	16	14	14	4	272	2	0	129
	Sparus_aurata	1	80	0	7	0	0	0	1	0	0	217	3	0	0	0	37	6	0	0	6	29
	Sciaena_umbra	6	3	0	0	0	4	1	3	0	273	0	7	0	0	0	9	1	0	0	0	6
	Serranus_cabrilla	3	4	13	0	0	0	1	8	101	0	0	2	28	3	2	0	0	1	5	1	22
	Anthias_anthias	90	7	2	23	38	0	0	348	6	2	7	33	54	2	13	4	0	0	0	0	121
	Serranus_scriba	0	1	0	0	0	0	105	0	2	0	0	1	0	1	0	0	0	0	6	0	6
	Epinephelus_marginatus	3	2	0	0	0	151	0	0	1	18	0	3	1	2	0	5	1	0	0	0	10
	Apogon_imberbis	8	0	0	0	115	0	0	48	1	0	0	1	1	0	0	0	0	0	0	0	31
	Sarpa_salpa	1	2	0	1966	0	0	0	1	1	1	3	13	5	11	0	33	4	0	0	6	122
	Mullus_surmuletus	0	0	78	1	0	0	0	0	0	0	0	0	3	4	0	0	0	0	0	0	1
	Diplodus_vulgaris	8	1640	0	10	0	2	4	5	0	1	6	2	2	0	1	93	0	0	0	9	254
	Chromis_chromis	7949	5	1	390	24	0	4	117	22	6	1	19	28	17	64	10	4	85	1	1	2479
	Chromis_chromis																					
	Diplodus_vulgaris																					
	Mullus_surmuletus																					
	Sarpa_salpa																					
	Apogon_imberbis																					
	Epinephelus_marginatus																					
	Serranus_scriba																					
	Anthias_anthias																					
	Serranus_cabrilla																					
	Sciaena_umbra																					
	Sparus_aurata																					
	Other																					
	Coris_julis																					
	Symphodus_tinca																					
	Spicara_smaris																					
	Diplodus_sargus																					
	Seriola																					
	Boops_boops																					
	Thalassoma_pavo																					
	Oblada_melanura																					
	background FP																					

Figure 3: The confusion matrix of the presented YOLOv7 model describing the relation between Predicted and True classes not in relation but in absolute values.

Three different examples of imagery and the corresponding detections of the YOLO model are presented. Original extracted frames without any preprocessing are presented

in figure 4-A, whilst 4-B depicts preprocessed images with the corresponding detections. Two images depict a rocky substrate with a variety of different species seen in the area (primarily *S. salpa* & *Sciaena umbra*) and the third image shows an artificial reef with a school of *Diplodus vulgaris* passing. Partial Figures in 4-C are framed as correctly detected species in green and missed detections in red.

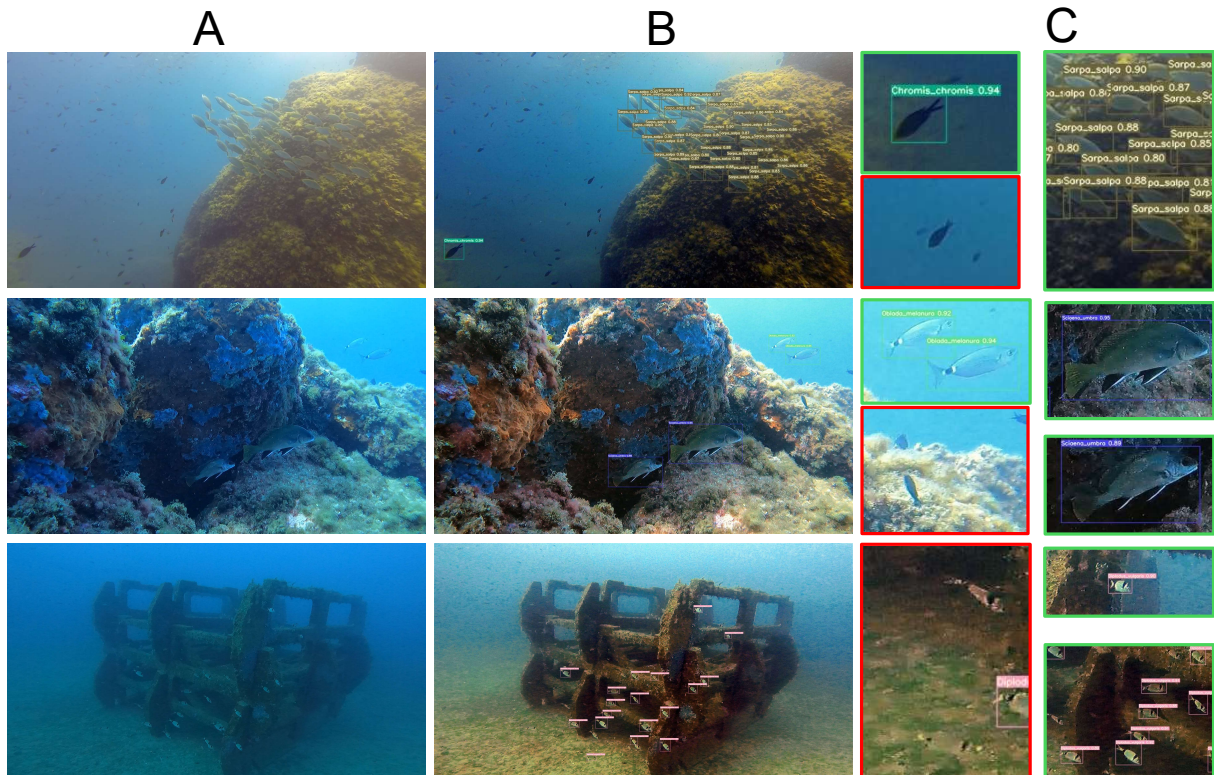


Figure 4: Example imagery of the dataset (A). Preprocessing and corresponding detections (B). Column C corresponds to good detections (framed in green) and missed detections (framed in red).

### 3.2 Binary presence/absence task evaluation

For the second experiment, the experimental dataset ( $DATA_E$ ) was used from the field campaign the year after  $DATA_T$  was recorded. The model detections on the experimental dataset were formatted into a presence/absence table that showed the existence or non-existence of each species in a transect and was compared to the manual video counting method. Figure 5 shows the relation between sensitivity and 1-specificity of the model in predicting aforementioned presence/absence. This graph shows an area under the curve (AUC) of 0.93 which is sufficient in capturing the models performance towards the presented presence/absence task.

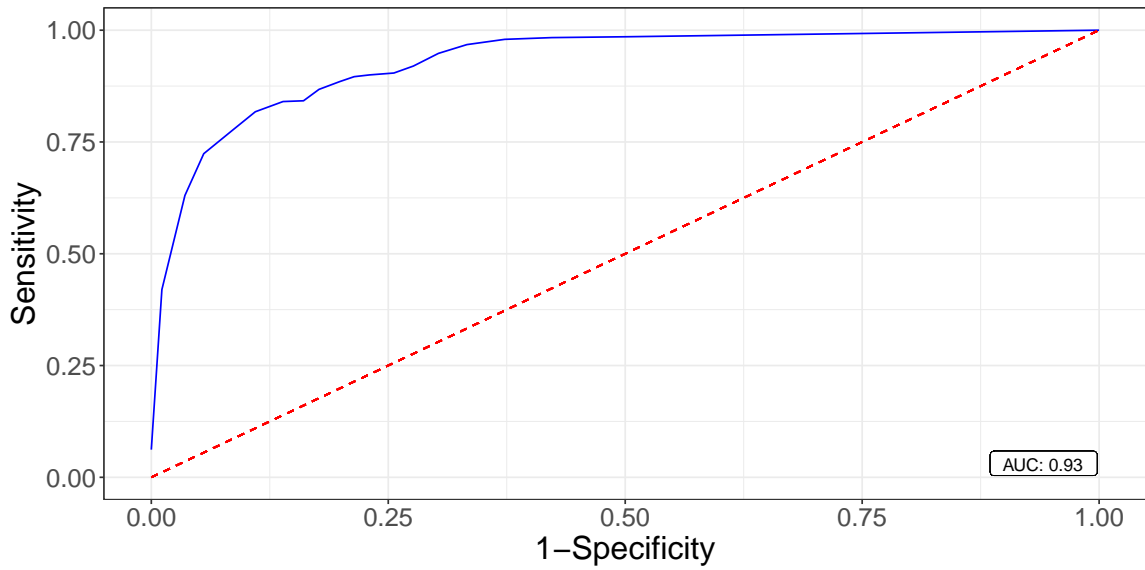


Figure 5: The sensitivity (blue line) is plotted against 1-specificity over 0.5 intervals of confidence thresholds allowing to extract the AUC indicating the performance of the model in predicting the species-transect-pairs overall. The AUC number is indicated on the right bottom corner and the dashed red line corresponds to a random classification.

The species-specific AUC ranges (Fig. 6) from 0.36 (*Serranus cabrilla*) to 0.99 (*Diplodus vulgaris*, *Seriola sp.*, *Thalassoma pavo* & *Apogon imberbis*) for the different species. An AUC of 0.00 (*S. umbra* & *Sparus aurata*) meaning that there was no observation of this species in the transect and are therefore ignored. Not all species curves behave the same way and some species are easier to detect in the transect when in comparison with harder to detect *M. surmuletus* or more elusive *S. cabrilla* species.

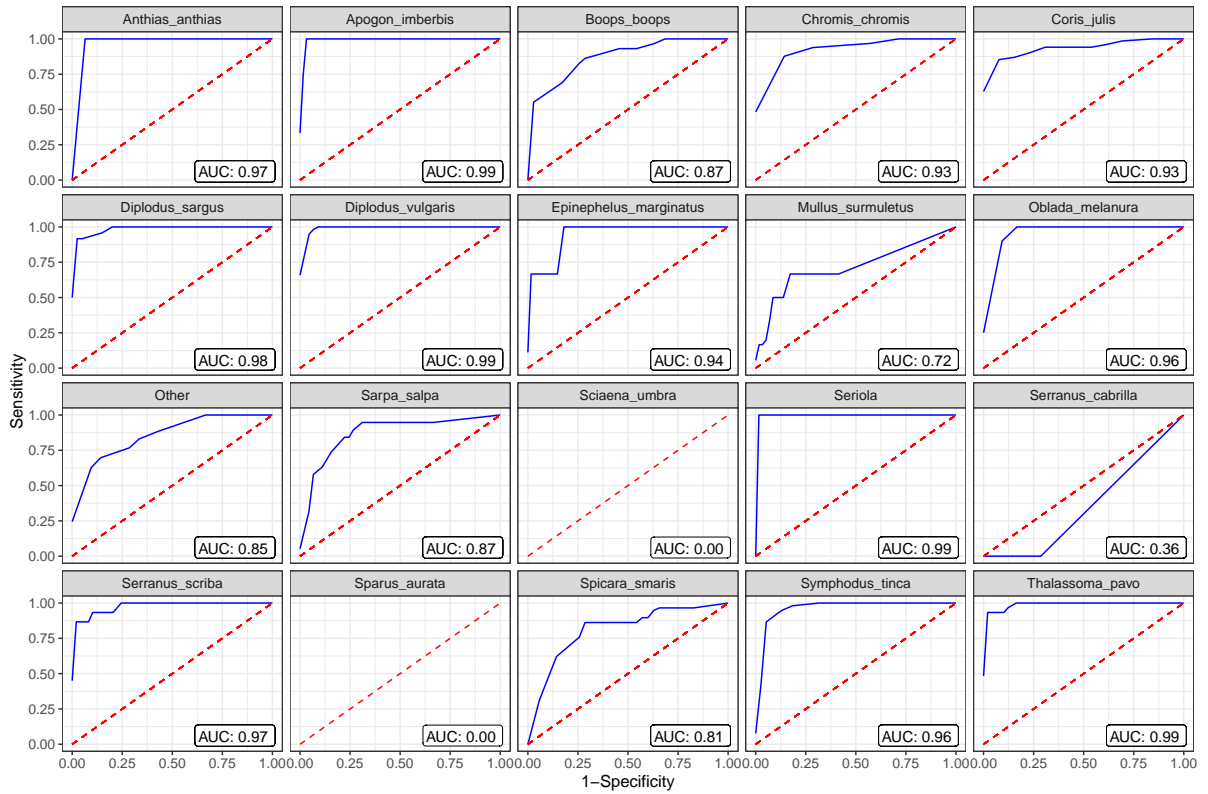


Figure 6: The sensitivity (blue line) is plotted against 1-specificity over 0.5 intervals of confidence thresholds allowing to extract the AUC indicating the performance of the model in predicting the species-transect pairs. Indicating how well the model is able to grasp the species richness of each transect for each of the species investigated. The AUC number is indicated on the right bottom corner and the dashed red line corresponds to a random classification.

The presence/absence per species per transect (referred to as species-transect pair) was assessed for different binary metrics (accuracy, sensitivity & specificity). These metrics are presented in relation to the confidence threshold (a value that represents the minimum level of certainty required for the model to classify an object) in Figure 7, which gives a better indication how the detection confidence has an impact on the corresponding model metric. Any detection with a confidence lower than the threshold will be disregarded and not incorporated into the further evaluation increasing potential FN but also decreasing FP. This threshold can and should be chosen according to the task presented and what the project managers requirements are. Purple dashed lines in Figure 7 indicate the different thresholds chosen for the transect evaluation and look at the impact of choosing these. The threshold 0.05 was chosen as loose to ensure all individuals are detected, 0.60 was chosen as balanced meant to balance the proportion of FP and FN (the crossing of sensitivity and specificity) and 0.80 was chosen as strict to be sure that mostly TP are included and FP are excluded in the analysis.

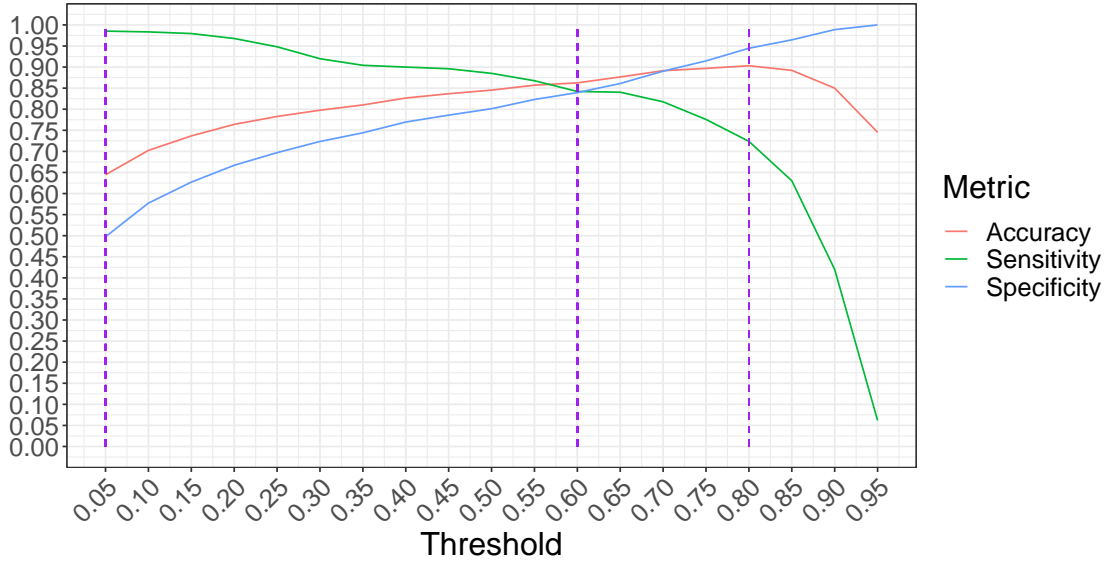


Figure 7: Three different metrics - accuracy (red), sensitivity (green) & specificity (blue) - are plotted against the confidence thresholds. This allows to evaluate the influence of choosing different thresholds for the transect evaluation. The purple dashed lines indicate the thresholds (loose, balanced & strict) chosen for the further study.

Depending on the threshold chosen, the results in Table 1 differ in the metrics mentioned previously. Accuracy values range from 0.64 to 0.90, sensitivity from 0.99 to 0.72 and specificity from 0.50 to 0.94 for the different thresholds. Accuracy and specificity are highest with the strict threshold while sensitivity is highest in the loose threshold. The stricter the threshold the higher the specificity, while the sensitivity increases with a more loose threshold. Cohens Kappa is highest for the strict threshold and TSS is highest for the more balanced dataset with values of 0.69 for Cohens Kappa and 0.68 for TSS respectively.

Table 1: Evaluation metrics per threshold analysis for all the transects. Bold numbers indicate the highest values for each metric.

Threshold	Accuracy	Sensitivity	Specificity	Cohens Kappa	TSS
Loose (0.05)	0.64	<b>0.99</b>	0.50	0.34	0.48
Balanced (0.60)	0.86	0.84	0.84	0.64	<b>0.68</b>
Strict (0.80)	<b>0.90</b>	0.72	<b>0.94</b>	<b>0.69</b>	0.67

The three different methods (manual video data, automated video data and diver data) were evaluated using color codes (Fig. 8) - red indicating a FN, yellow indicating a FP and orange indicating a TP. Transparent squares are TN. Squares marked with a cross are species that have been seen by the diver in the transect during the UVC. Transects 1-39 are over seagrass meadows (*Posidonia oceanica*) while transects 39 to 64 are over rocky substrates. The average species richness per site for the two DL-independent methods over rocky substrates are 7.88 (manual video count) and 9.12 (traditional UVC). Over seagrass meadows the average species richness is less with values of 3.89 and 5.05, respectively.



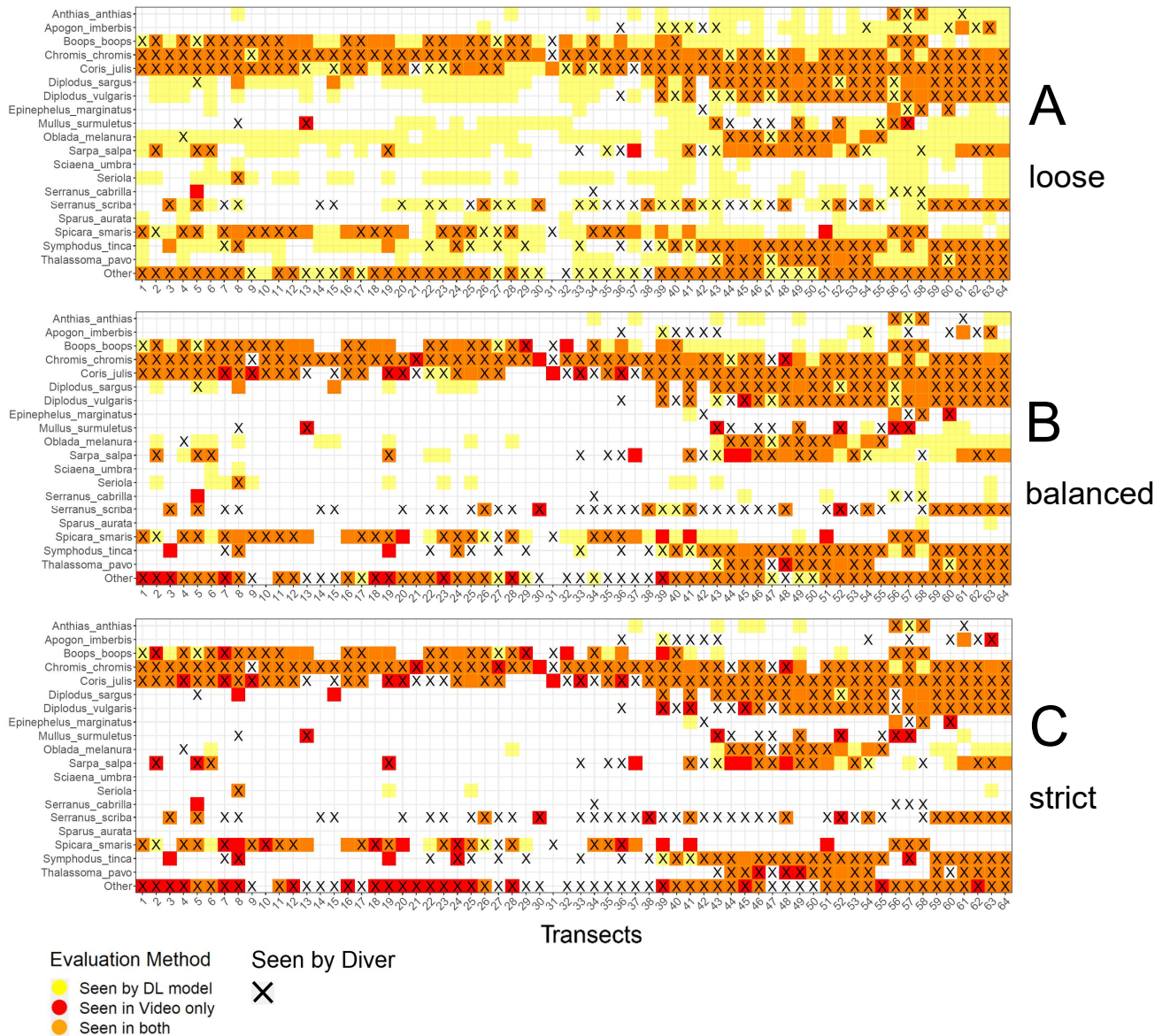


Figure 8: Transect evaluation with the different thresholds: A, B and C are loose, balanced and strict thresholds respectively. Successful detections are indicated in orange, red are missed species (FN) and in yellow falsely detected species (FP). Cross-marked squares show presence of the corresponding fish species in the transect as seen by the diver doing the UVC.

In the threshold analysis, the automatically generated DL-derived presence/absence table were compared. For the loose threshold (Fig. 8 - A) the transect evaluation includes a total of 450 FP (yellow colour) while having 5 FN (red colour). In total 348 species-transect pairs were correctly predicted (TP) and 477 were correctly not predicted (TN). Even though there is a high number of FP, the species *M. surmuletus* was missed in 2 out of the 6 appearing videos while appearing 24 time as an FP and is therefor the worst performing species. This result is in alignment with the confusion matrix of the model evaluation (Fig. 3) with 121 detections of *M. surmuletus* missed. The different substrates are not distinguishable in the loose threshold case. The predicted average species richness

over rocky substrates is 15.5 and 10.39 for seagrass meadows. The species richness is overestimated with 7.62 species for rocky substrates and 6.5 for the seagrass meadows.

The balanced threshold model configuration (Fig. 8 - B) gives 135, 41, 312 and 792, FP, FN, TP and TN, respectively. Out of 1280 possible predictions, 1104 were predicted correctly by the DL model. *M. surmuletus* has been seen by the DL model in 1 of the 6 videos this species is appearing in whilst having 3 FP. Regarding the predicted species richness over rocky substrates the value is 10.54. The species richness over seagrass meadows corresponds to 4.68. There is a clear difference observable - qualitatively and quantitatively - between the two substrate types. Differences in species richness in automation and manual count are 2.66 for the rocky substrates and 0.79 for the seagrass meadows.

The FP, FN, TP and TN for the strict threshold (Fig. 8 - C) were valued at 45, 79, 274 and 882, respectively. A total of 1156 species-transect-pairs have been correctly predicted. *M. surmuletus* has been seen by the DL model in 1 of the 6 appearing videos for the strict threshold and performing the best in this strict case with 1 FP. For the strict case, the predicted species richness over rocky substrates and seagrass meadows are 8.04 and 3.06, respectively. The difference in species richness is less visible with this threshold with species richness differences of -0.83 species over seagrass substrates and 0.16 species richness difference for rocky substrates. The strict threshold is the only threshold that differs in less than 1 over both substrates but underestimates the average seagrass richness when compared to the manual video count.

### 3.3 Method comparison

The multi-dimensional, quasi-proportional Venn diagrams from the R package nVennR (Pérez-Silva *et al.*, 2018) display a qualitative insight on how the sites are evaluated with all the techniques and how large the agreement between them is. Depending on availability of video recordings and diver data, 3 to 6 transects were concatenated *per* site meaning the transects were fused to create the site-specific fish diversity. For the DL part the strict confidence threshold was chosen since the accuracy (1156 species-transect pairs were correctly predicted) was high whilst having substantial Cohens Kappa, excellent TSS values and better species richness prediction.

From the models point of view, a yellow & green-yellow colour indicates FP predictions since these species were not seen in the video but were detected by the DL model. Distinct red & red-green color specifies FN since there was a presence in the video but not detected from the model. The yellow-red & all three color overlap show TP. A distinct green coloration shows a species that has only been seen by the diver during the UVC and will not flow into the evaluation of the model but will be used to give an idea on how the different methods add to a site diversity. A total of 26 FN and 16 FP were predicted while

94 TP and 144 TN form a total of 280 species-site pairs.

An overall agreement can be observed between the methods for the location Cap Roux (Fig. 9) concerning the species richness. Most species per site have been counted with each of the methods evaluated. The FP (yellow & green-yellow) are ranging from 0 to 2 with a median of 2 for the DL model (Fig. 9 - b, d & f-i). The green-yellow could show human error when evaluating the videos but were double checked to be correctly classified as FP. Missed species in the video (red) were observed on 6 sites (Fig. 9 - a-c, f & i) while on 6 sites (Fig. 9 - a-c, g & h) the diver has seen fish species that were not observed in the video (green). In Figure 9 - d the DL model was agreeing with the manual video count while the diver missed this species (yellow-red). An overlap between the diver and the manual video count is presented in 6 sites (Fig. 9 - a, c, d, f, g & h).

The overlapping average species richness for seagrass meadows at all depths is 4.6 while the species richness of rocky habitats at all depths is 9.5.

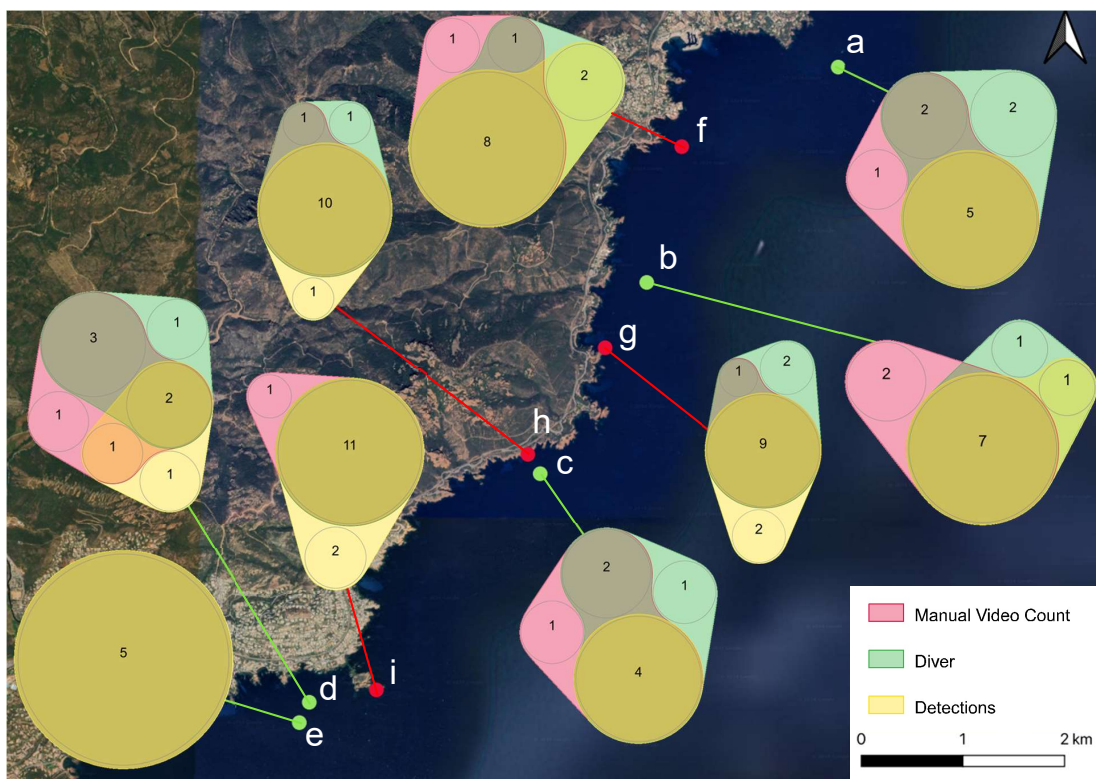


Figure 9: Venn diagrams at the site coordinates (green sites (a-e) are seagrass meadows, red sites (f-i) are rocky substrates) showing the overlap of methods in different sites at the location **Cap Roux** (Image credit Google Maps 2024). Overlapping colors indicate a common species pool whilst distinct coloration indicates method-specific detection. Yellow & yellow-green colour are indicating FP and red & red-green colour indicates FN. Green coloration means the video has not sufficiently covered all species in the transect and was only seen by the diver.

For the second location in Corniche Varoise (Fig. 10) the video transects do not



contribute identifying new species that are unseen by the diver. However, on 4 of the 5 (Fig. 10 - b-e) sites species have only been seen by the diver with missed species ranging from 1 to 5 with a median of 2. There are DL model-derived FP at all sites ranging from 1 to 3 with a median of 1. Compared to Cap Roux an increase in average FP and an overall decrease in video-contributed species is observable.

Overall the seagrass meadow sites have a species richness of 6.5 and the rocky substrate sites have a species richness of 9.5.

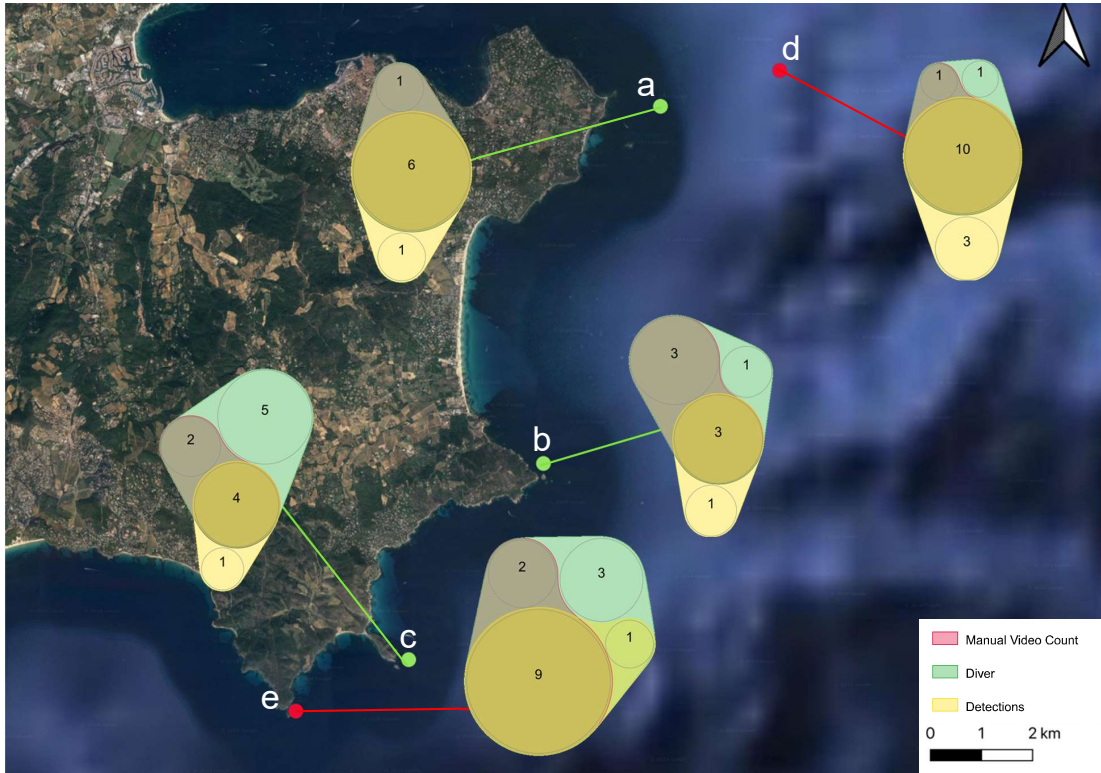


Figure 10: Venn diagrams at the site coordinates (green sites (a-c) are seagrass meadows, red sites (d-e) are rocky substrates) of the site showing the overlap of methods in different sites at the location **Corniche Varoise** (Image credit Google Maps 2024). Overlapping colors indicate a common species pool whilst distinct coloration indicates method-specific detection. Yellow & yellow-green colour are indicating FP and red & red-green colour indicates FN. Green coloration means the video has not sufficiently covered all species in the transect and was only seen by the diver.

## 4 Discussion

Overall this study shows the successful incorporation of a DL model into the validation and evaluation process of assessing biodiversity at 2 different ecologically important sites of the French Riviera. The proposed model is able to detect the majority of the presence or absence of 19 different fish species in each transect. This makes it a valid addition to an holistic approach and increases robustness of justification for MPA implementation. The

automation of the video analysis will decrease labour time and increase the time spent on interpretation and human-dependent decision-making process increasing the efficiency of the whole procedure. This highlights the importance of utilising DL tools in applied scenarios and by simplifying the task of transect-based, diver operated videos (DOV) evaluation.

The proposed deep learning model provides evidence for a sufficient detection capability on DOV. With an mAP value of 0.56 this study allows a benchmark for more work on similar data in other regions which is scarce to this day with the exceptions. For example Connolly *et al.*, 2022 uses videos from a ROV which replaces the diver and shows good results in detecting two species on the Australian coast. Reporting species precision values of 0.28 to 0.96 and species recall values of 0.46 to 0.80 which, in comparison to the proposed model precision (0.13 to 0.91) and recall (0.23 to 0.75), is in agreement but is differently challenged with the task of detecting 2 species in comparison to the 20 species proposed in this study. In comparison to stationary camera studies (Jalal *et al.*, 2020; Knausgård *et al.*, 2022) both aforementioned studies on non-stationary cameras underperform due to the difference in data collection and challenge in the movement of the camera. Other studies have a more constrained environment and achieve very inspiring results in for example re-identifying individuals (Olsen *et al.*, 2023) that could be very interesting for a future collaborative step when concerning *E. marginatus* due to its big size and importance as a top predator and protected species. Marine biologist rely on fast and efficient information on all of their data.

With a different task, the proposed pipeline shows insights into the transect data without prior knowledge and creating therefor a simplified, more holistic approach. Most of the species reached an AUC value of over 0.80 which seems sufficient in detecting the presence/absence in a transect. However, limitations are still prevalent since the detection rates of certain species are insufficient to confidently predict the presence of species in the transects - for example *M. surmuletus* with an AUC value 0.72 which is an economically important species in this region. While *M. surmuletus* shows insufficient AUC values, *i.e.* *A. imberbis* shows great coverage in the transect analysis with an AUC value of 0.99. This AUC value makes us believe that the model is predicting the species *A. imberbis* perfectly which is not evident. The species *A. imberbis* only appeared in 2 transects and was detected in this transect at most confidence thresholds which gives a sensitivity of 1. Lower thresholds hold some FP and therefor, decrease the specificity but since the sensitivity is almost perfect in most of the threshold-dependent cases the AUC is close to 1. In our case more species fall into this low-occurrence category and should not be considered with the AUC evaluation - *Anthias anthias*, *Epinephelus marginatus*, *Seriola sp.* and *S. cabrilla*. Having an imbalance between TP and TN can lead to misinterpretation of a model and

other metrics should be consolidated (Lobo *et al.*, 2008). For example the Cohens Kappa gives a better insight on the overall inter-reliability between two measurement methods and to what degree they agree with each other but suffers from prevalence. The TSS metric on the other hand is free of this prevalence and should be considered as well to further justify metric values of the model. In our case the highest value for the Cohens Kappa is 0.69 which is considered substantial (Landis and Koch, 1977) and the highest True Skill Statistic (TSS) value of 0.68 is considered excellent (Komac *et al.*, 2016). These two metrics show that there is a high degree of agreement between the manual & automated method which shows evidence that the manual video labour could be replaced with a DL algorithm.

With prior mentioned information on the confidence thresholds, the method-specific species evaluation was conducted with the strict threshold. It showed the highest accuracy (0.90) as well as the best Cohens Kappa (0.69) and second best TSS (0.67). With this in mind, there will be an increase in FN but also a substantial decrease in FP which is important to correctly estimate species richness on transects. It is important to evaluate two cases, firstly only UVC and video transects and secondly how the addition of a DL model changes the overall species richness. For the first case, at the sites of Corniche Varoise location (Fig. 10) there are only added species from the traditional UVC which highlights the importance of expert knowledge (Cappo *et al.*, 2003). However, the manual video count added species at Cap Roux (red colored) and there is a good overlap with additions to the species richness from both methods (red-green colored in Fig. 9). This further shows evidence of the importance of doing a more holistic approach (Colton and Swearer, 2010) and the need to adapt the video collection protocol to bear more resemblance to the traditional UVC data collection. Reasons for these differences could be on one hand that the UVC is more strict in terms of field of view and it limits the diver on purpose to a given surface area of 125m<sup>2</sup> (Grorud-Colvert *et al.*, 2021) explaining the additions of species by video methods at Cap Roux. While on the other hand, the divers head is free in its movement and only limited by theoretical area given by the protocol, the cameras point of view is always limited. This means that a diver can and does eliminate these spatial limitation of a camera that could explain the addition of species at at the location of Corniche Varoise. These differences can have a small but significant impact on a species being sighted in a transect or not - by both methods. This difference is more prevalent for rarer species since they can be missed easier and not reappear later in the transect. Therefore, adaptations to these differences should be carried out and could be for example using multiple cameras, one facing straight ahead and the other one looking down at the seabed to cover both semi-pelagic and benthic species closely observing what the diver observes. Another protocol modification could be to increase the area covered of DOV transects to enable a fair comparison with the UVC. Further research needs to be

conducted to find the correct protocol parameters for fair comparison.

For the second case we look at the performance of the automation of the video transect evaluation. While other studies have shown to successfully evaluate stationary or baited videos (Villon *et al.*, 2018; Jalal *et al.*, 2020; Muksit *et al.*, 2022), the proposed study utilised transect-derived and depth-variable DOV, which - as mentioned in Villon *et al.*, 2018 & Connolly *et al.*, 2022 - is another important use-case of artificial intelligence in marine conservation and should be the subject of increased focus in the future. When adding a high throughput DL model on transect videos instead of manual labour, a good level of quality is kept overall while minimising analysis time. A total of 85% of the species-site pairs have been correctly detected which highlights the feasibility of this approach and eliminating the divers bias which can correspond to up to 25 % of a local species diversity. The proposed DL model is able to grasp the reality and is an indication that in the future, this pipeline helps to gather more bias-free data for MPA managers to streamline the process of implementing and evaluating MPA, making the decision-making process faster and more efficient.

Besides the data on biodiversity and its evolution over time, different ecological information can be extracted from this simplified task that are interesting for aforementioned establishments of protected areas or no take zones. While seagrass meadows are known for nursery and living habitats (Franco *et al.*, 2006; Boudouresque *et al.*, 2021), several studies have shown that the species richness is higher over rocky substrate rather than *P. oceanica* (Cheminée *et al.*, 2021). This difference was also observed in this study with the focalised species list that the DL model was trained on. The protected species *E. marginatus* was seen by the DL model only over rocky substrate and occurred twice as much inside of the protected area than outside. This gives an indication that an MPA can be an important safe haven for top predator fishes in their hunting and reproduction grounds which is in alignment with previous studies (Guidetti *et al.*, 2014). In the case of the two *Diplodus* *ssp.* the analysis showed that these two economically important species are represented mostly over rocky substrate which further highlights the importance of this habitat and how artificialisation of the rocky shoreline can damage different species more than others. These results are preliminary and need to be enjoyed with caution since the data is limited to two sites and one time frame. However, this pipeline will be applied to more and more frequent data collections and will generate information on different ecosystems, sites and time frames creating a standardised procedure for the future.

## 4.1 Conclusion

This study showed that no method of evaluation (video data nor diver data) is perfect but a combination of them would achieve the highest robustness over a site. A combinative

approach is in agreement with non-automated studies which looked at different methods by hand (Willis and Babcock, 2000; Goetze *et al.*, 2015). In the case of diver limitations (Cappo *et al.*, 2003), the video recording protocol needs to be adapted to create data that is comparable between methods and create a strong, robust and reproducible protocol. Furthermore showed in this study, deep learning models can assist the researcher to disconnect the video analysis from any expert bias and save time and labour due to the automation of the procedure. With the decrease in labour time, more frequent data collection missions can be targeted increasing the temporal and spatial coverage of species distribution. This proof of concept shows not only a new high-throughput analysis of diver operated videos (DOV) but also the possibility to evaluate videos recorded by remote operated vehicles (ROV) which allows data to be collected in scenarios where there is a meteorological, depth-related or logistical limitation.

## Acknowledgements

This work was only made possible thanks to the collaboration with the projects RECIF (Réseau d’Evaluation des Cantonnements et ZSC en Interface Fonctionnelle) and FEAMPA (Fonds européen pour les affaires maritimes, la pêche et l’aquaculture) and the divers involved who provided the diver- and video data from the corresponding field campaigns. The authors are grateful to the OPAL infrastructure from Université Côte d’Azur for providing resources and support. This project was funded through the UCAJEDI Investments in the Future project managed by the National Research Agency (ANR) with the reference number ANR-15-IDEX-01.

## References

- AIMS, University of Western Australia, and Curtin University (2019). *OzFish Dataset - Machine learning dataset for Baited Remote Underwater Video Stations*. URL: <https://apps.aims.gov.au/metadata/view/38c829d4-6b6d-44a1-9476-f9b0955ce0b8>.
- Allken, V. *et al.* (2021a). “A deep learning-based method to identify and count pelagic and mesopelagic fishes from trawl camera images”. In: *ICES Journal of Marine Science* 78, pp. 3780–3792. DOI: 10.1093/icesjms/fsab227.
- (2021b). “A real-world dataset and data simulation algorithm for automated fish species identification”. In: *Geoscience Data Journal* 8, pp. 199–209. DOI: 10.21335/NMDC-551736490.
- Allouche, O., A. Tsoar, and R. Kadmon (2006). “Assessing the accuracy of species distribution models: prevalence, Kappa and the True Skill Statistic (TSS)”. In: *Journal of Applied Ecology* 43, pp. 1223–1232. DOI: 10.1111/j.1365-2664.2006.01214.x.
- Boudouresque, C.-F. *et al.* (2021). “Restoration of Seagrass Meadows in the Mediterranean Sea: A Critical Review of Effectiveness and Ethical Issues”. In: *Water* 13, p. 1034. DOI: 10.3390/w13081034.
- Caldwell, Z. R. *et al.* (2016). “Reef fish survey techniques: assessing the potential for standardizing methodologies”. In: *PloS one* 11, e0153066. DOI: 10.1371/journal.pone.0153066.
- Cappo, M. *et al.* (2003). “Aquatic Protected Areas - what works best and how do we know”. In: chap. Potential of video techniques to monitor diversity, abundance and size of fish in studies of Marine Protected Areas, pp. 455–464. URL: [https://books.google.fr/books?id=9\\_NEAAAAYAAJ](https://books.google.fr/books?id=9_NEAAAAYAAJ).
- Carranza, M. *et al.* (2019). “Urban expansion depletes cultural ecosystem services: an insight into a Mediterranean coastline”. In: *Rendiconti Lincei. Scienze Fisiche e Naturali* 31, pp. 103–111. DOI: 10.1007/s12210-019-00866-w.
- Cheminée, A. *et al.* (2021). “All shallow coastal habitats matter as nurseries for Mediterranean juvenile fish”. In: *Scientific Reports* 11, p. 14631. DOI: 10.1038/s41598-021-93557-2.
- Colton, M. and S. Swearer (2010). “A comparison of two survey methods: Differences between underwater visual census and baited remote underwater video”. In: *Marine Ecology-Progress Series (MEPS)* 400, pp. 19–36. DOI: 10.3354/meps08377.
- Connolly, R. M. *et al.* (2022). “Fish surveys on the move: Adapting automated fish detection and classification frameworks for videos on a remotely operated vehicle in shallow marine waters”. In: *Frontiers in Marine Science* 9. DOI: 10.3389/fmars.2022.918504.
- Demirel, N., M. Zengin, and A. Ulman (2020). “First Large-Scale Eastern Mediterranean and Black Sea Stock Assessment Reveals a Dramatic Decline”. In: *Frontiers in Marine Science* 7. DOI: 10.3389/fmars.2020.00103.

- Doney, S. C. *et al.* (2012). “Climate Change Impacts on Marine Ecosystems”. In: *Annual Review of Marine Science* 4, pp. 11–37. DOI: 10.1146/annurev-marine-041911-111611.
- Franco, A. *et al.* (2006). “Use of shallow water habitats by fish assemblages in a Mediterranean coastal lagoon”. In: *Estuarine, Coastal and Shelf Science* 66, pp. 67–83. DOI: 10.1016/j.ecss.2005.07.020.
- Garcia-D’Urso, N. *et al.* (2022). “The DeepFish computer vision dataset for fish instance segmentation, classification, and size estimation”. In: *Scientific Data* 9, p. 287. DOI: 10.1038/s41597-022-01416-0.
- Gissi, E. *et al.* (2021). “A review of the combined effects of climate change and other local human stressors on the marine environment”. In: *Science of The Total Environment* 755, p. 142564. DOI: 10.1016/j.scitotenv.2020.142564.
- Goetze, J.S. *et al.* (2015). “Diver operated video most accurately detects the impacts of fishing within periodically harvested closures”. In: *Journal of Experimental Marine Biology and Ecology* 462, pp. 74–82. DOI: 10.1016/j.jembe.2014.10.004.
- Grorud-Colvert, K. *et al.* (2021). “The MPA Guide: A framework to achieve global goals for the ocean”. In: *Science* 373, eabf0861. DOI: 10.1126/science.abf0861.
- Guidetti, P. *et al.* (2014). “Large-Scale Assessment of Mediterranean Marine Protected Areas Effects on Fish Assemblages”. In: *PloS one* 9, e91841. DOI: 10.1371/journal.pone.0091841.
- Harmelin-Vivien, M. L. *et al.* (1985). “Evaluation visuelle des peuplements et populations de poissons méthodes et problèmes”. In: *Revue d’Écologie (La Terre et La Vie)* 40, pp. 467–539. DOI: 10.3406/revec.1985.5297.
- Jalal, A. *et al.* (2020). “Fish detection and species classification in underwater environments using deep learning with temporal information”. In: *Ecological Informatics* 57, p. 101088. DOI: 10.1016/j.ecoinf.2020.101088.
- Knausgård, K. *et al.* (2022). “Temperate fish detection and classification: a deep learning based approach”. In: *Applied Intelligence* 52, pp. 6988–7001. DOI: 10.1007/s10489-020-02154-9.
- Komac, B. *et al.* (2016). “Modelization of the Current and Future Habitat Suitability of *Rhododendron ferrugineum* Using Potential Snow Accumulation”. In: *PloS one* 11, e0147324. DOI: 10.1371/journal.pone.0147324.
- Kulbicki, M. and S. Sarramégna (1999). “Comparison of density estimates derived from strip transect and distance sampling for underwater visual censuses: a case study of Chaetodontidae and Pomacanthidae”. In: *Aquatic Living Resources* 12, pp. 315–325. DOI: 10.1016/S0990-7440(99)00116-3.
- Landis, J. R. and G. G. Koch (1977). “The Measurement of Observer Agreement for Categorical Data”. In: *Biometrics* 33, pp. 159–174. DOI: 10.2307/2529310.

- Lobo, J. M., A. Jiménez-Valverde, and R. Real (2008). “AUC: a misleading measure of the performance of predictive distribution models”. In: *Global Ecology and Biogeography* 17, pp. 145–151. DOI: 10.1111/j.1466-8238.2007.00358.x.
- Malde, K. *et al.* (2020). “Machine intelligence and the data-driven future of marine science”. In: *ICES Journal of Marine Science* 77, pp. 1274–1285. DOI: 10.1093/icesjms/fsz057.
- Mcclanahan, T. *et al.* (2007). “Influence of instantaneous variation on estimates of coral reef fish populations and communities”. In: *Marine Ecology Progress Series* 340, pp. 221–234. DOI: 10.3354/meps340221.
- Mclean, D. *et al.* (2005). “A comparison of temperate reef fish assemblages recorded by three underwater stereo-video techniques”. In: *Marine Biology* 148, pp. 415–425. DOI: 10.1007/s00227-005-0090-6.
- Mejjad, N., A. Rossi, and A. B. Pavel (2022). “The coastal tourism industry in the Mediterranean: A critical review of the socio-economic and environmental pressures impacts”. In: *Tourism Management Perspectives* 44, p. 101007. DOI: 10.1016/j.tmp.2022.101007.
- Mohamed, H. E.-D. *et al.* (2020). “MSR-YOLO: Method to Enhance Fish Detection and Tracking in Fish Farms”. In: *Procedia Computer Science* 170, pp. 539–546. DOI: 10.1016/j.procs.2020.03.123.
- Muksit, A. A. *et al.* (2022). “YOLO-Fish: A robust fish detection model to detect fish in realistic underwater environment”. In: *Ecological Informatics* 72, p. 101847. DOI: 10.1016/j.ecoinf.2022.101847.
- Myers, R. A., J. A. Hutchings, and N. J. Barrowman (1997). “Why do Fish Stocks Collapse? The Example of Cod in Atlantic Canada”. In: *Ecological Applications* 7, pp. 91–106. DOI: 10.2307/2269409.
- Olsen, Ø. L. *et al.* (2023). “A contrastive learning approach for individual re-identification in a wild fish population”. In: DOI: 10.48550/arXiv.2301.00596.
- Park, J.-H. and C. Kang (2020). “A Study on Enhancement of Fish Recognition Using Cumulative Mean of YOLO Network in Underwater Video Images”. In: *Journal of Marine Science and Engineering* 8. DOI: 10.3390/jmse8110952.
- Pérez-Silva, J. G., M. Araujo-Voces, and V. Quesada (2018). “nVenn: generalized, quasi-proportional Venn and Euler diagrams”. In: *Bioinformatics* 34, pp. 2322–2324. DOI: 10.1093/bioinformatics/bty109.
- Priyankan, K. and T. G. I. Fernando (2021). “Mobile Application to Identify Fish Species Using YOLO and Convolutional Neural Networks”. In: *Proceedings of International Conference on Sustainable Expert Systems*. Springer Singapore, pp. 303–317. DOI: 10.1007/978-981-33-4355-9\_24.
- Rubbens, P. *et al.* (2023). “Machine learning in marine ecology: an overview of techniques and applications”. In: *ICES Journal of Marine Science* 80, pp. 1829–1853. DOI: 10.1093/icesjms/fsad100.



- Smale, D. *et al.* (2019). “Marine heatwaves threaten global biodiversity and the provision of ecosystem services”. In: *Nature Climate Change* 9, pp. 306–312. DOI: 10.1038/s41558-019-0412-1.
- Spampinato, C. *et al.* (2016). “Fine-grained object recognition in underwater visual data”. In: *Multimedia Tools and Applications* 75, pp. 1701–1720. DOI: 10.1007/s11042-015-2601-x.
- Vabø, R. *et al.* (2021). “Automatic interpretation of salmon scales using deep learning”. In: *Ecological Informatics* 63, p. 101322. DOI: 10.1016/j.ecoinf.2021.101322.
- Vasilakopoulos, P., C. D. Maravelias, and G. Tserpes (2014). “The alarming decline of Mediterranean fish stocks”. In: *Current Biology* 24, pp. 1643–1648. DOI: 10.1016/j.cub.2014.05.070.
- Villon, S. *et al.* (2018). “A Deep learning method for accurate and fast identification of coral reef fishes in underwater images”. In: *Ecological Informatics* 48, pp. 238–244. DOI: 10.1016/j.ecoinf.2018.09.007.
- Wang, C.-Y., A. Bochkovskiy, and H.-Y M. Liao (2023). “YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors”. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 7464–7475. DOI: 10.48550/arXiv.2207.02696.
- Wang, Y. *et al.* (2021). “UIEC<sup>2</sup>-Net: CNN-based underwater image enhancement using two color space”. In: *Signal Processing: Image Communication* 96, p. 116250. DOI: 10.1016/j.image.2021.116250.
- Willis, T. and R. Babcock (2000). “A baited underwater video system for the determination of relative density of carnivorous reef fish”. In: *Marine and Freshwater Research* 51, pp. 755–763. DOI: 10.1071/MF00010.
- Wu, C. M. *et al.* (2022). “Underwater trash detection algorithm based on improved YOLOv5s”. In: *Journal of Real-Time Image Processing* 19, pp. 911–920. DOI: 10.1007/s11554-022-01232-0.
- Xu, W. and S. Matzner (2018). “Underwater Fish Detection Using Deep Learning for Water Power Applications”. In: *5th Annual Conf. on Computational Science & Computational Intelligence (CSCI’18)*. DOI: 10.1109/CSCI46756.2018.00067.

## Supplementary Materials

Table S1: List of site specific meta data information in the French Riviera. At each site a different number of transects were conducted. They were randomly chosen over the correct substrates to have a random stratified sampling. The abbreviations I & O represent if the site is inside or outside of an MPA. Deep indicates depths below 15m and shallow indicate depths above 15m.

<b>Location</b>	<b>Substrate</b>	<b>Depth</b>	<b>Transects</b>	<b>I / O</b>	<b>Latitude</b>	<b>Longitude</b>
Cap Roux	posidonia	deep	1-6	I	43°27.753'N	6°55.728'E
Cap Roux	posidonia	deep	7-9	O	43°28.892'N	6°57.123'E
Cap Roux	posidonia	deep	10-12	O	43°25.424'N	6°53.193'E
Cap Roux	posidonia	deep	13-18	O	43°26.743'N	6°54.950'E
Cap Roux	posidonia	shallow	19-21	O	43°25.531'N	6°53.266'E
Corniche varoise	posidonia	shallow	22-27	I	43°15.899'N	6°42.642'E
Corniche varoise	posidonia	shallow	28-32	I	43°12.142'N	6°40.950'E
Corniche varoise	posidonia	shallow	33-38	I	43°10.084'N	6°39.003'E
Cap Roux	rocky	deep	39-42	O	43°28.472'N	6°55.982'E
Cap Roux	rocky	shallow	43-47	I	43°27.409'N	6°55.423'E
Cap Roux	rocky	shallow	48-51	O	43°25.597'N	6°53.756'E
Cap Roux	rocky	shallow	52-55	O	43°26.845'N	6°54.860'E
Corniche varoise	rocky	deep	56-58	I	43°16.025'N	6°41.844'E
Corniche varoise	rocky	deep	59-64	I	43°09.542'N	6°37.377'E

Table S2: List of taxonomic categories included in the 'Other' class of the model. The taxonomic level is as high as possible.

**Taxonomic Category**

*Antherina sp.*

*Blennidae sp.*

*Conger conger*

*Dicentrarchus labrax*

*Diplodus anularis*

*Gobidae sp.*

*Labrus merula*

*Labrus mixtus*

*Labrus viridis*

*Mugilidae sp.*

*Muraena helena*

*Phycis phycis*

*Scorpaena scrofa*

*Sphyræna sphyraena*

*Spicara maena*

*Spondylisoma cantharus*

*Symphodus mediterraneus*

*Symphodus melanocercus*

*Symphodus roissali*

*Symphodus rostratus*

*Trachurus sp.*

Table S3: The different classes and their corresponding label counts in the training & test set.

<b>Class</b>	<b>Count Train</b>	<b>Count Test</b>
<i>Chromis_chromis</i>	8'925	10'863
<i>Diplodus_vulgaris</i>	2'995	2'479
<i>Coris_julis</i>	3'174	2'473
<i>Symphodus_tinca</i>	2'426	785
<i>Spicara_smaris</i>	1'857	1'073
<i>Diplodus_sargus</i>	1'854	1'385
<i>Seriola</i>	383	162
<i>Boops_boops</i>	1'524	1'739
<i>Thalassoma_pavo</i>	447	235
<i>Oblada_melanura</i>	972	425
<i>Mullus_surmuletus</i>	378	301
<i>Sarpa_salpa</i>	3'934	3'266
<i>Apogon_imberbis</i>	1'280	270
<i>Epinephelus_marginatus</i>	749	212
<i>Serranus_scriba</i>	344	209
<i>Anthias_anthias</i>	4'030	821
<i>Serranus_cabrilla</i>	242	237
<i>Sciaena_umbra</i>	309	414
<i>Sparus_aurata</i>	391	332
<i>Other</i>	4'165	513
<b>Total</b>	<b>40'379</b>	<b>28'194</b>

Table S4: The different species-specific precision and recall values at 50% confidence interval. Bold numbers indicate the highest number and italic numbers indicate the lowest numbers of each precision and recall.

<b>Class</b>	<b>Precision</b>	<b>Recall</b>
<i>Sarpa_salpa</i>	<b>0.91</b>	0.54
<i>Mullus_surmuletus</i>	0.90	<i>0.23</i>
<i>Sciaena_umbra</i>	0.87	0.63
<i>Serranus_scriba</i>	0.86	0.48
<i>Diplodus_vulgaris</i>	0.81	0.61
<i>Epinephelus_marginatus</i>	0.77	0.70
<i>Coris_julis</i>	0.73	<b>0.75</b>
<i>Chromis_chromis</i>	0.71	0.73
<i>Thalassoma_pavo</i>	0.66	0.41
<i>Boops_boops</i>	0.60	0.61
<i>Apogon_imberbis</i>	0.56	0.41
<i>Sparus_aurata</i>	0.56	0.61
<i>Symphodus_tinca</i>	0.53	0.53
<i>Serranus_cabrilla</i>	0.52	0.40
<i>Diplodus_sargus</i>	0.50	0.63
<i>Spicara_smaris</i>	0.48	0.56
<i>Anthias_anthias</i>	0.46	0.40
<i>Oblada_melanura</i>	0.42	0.43
<i>Seriola</i>	0.41	0.70
<i>Other</i>	<i>0.13</i>	0.26
<b>Total</b>	<b>0.66</b>	<b>0.65</b>