

Save the Whales: Body Keypoints Detection for Automated Annotation

Suhaas Kiran DG
 University of Massachusetts
 Amherst, MA
 sdg@umass.edu

Siddarth Suresh
 University of Massachusetts
 Amherst, MA
 siddarthsure@umass.edu

Kedarnath Chimmad
 University of Massachusetts
 Amherst, MA
 kchimmad@umass.edu

Abstract

001 *In the realm of wildlife research, aerial photogrammetry*
 002 *emerges as a non-intrusive means to gather data on*
 003 *wild animals. However, the vast volume of data generated*
 004 *necessitates efficient processing methods. Specifically, an-*
 005 *notating key points on whale bodies for metrics like body*
 006 *mass proves to be time-consuming. To address this, we pro-*
 007 *pose a keypoint detection model based on the segmentation*
 008 *backbone and regression head. Using computer vision tech-*
 009 *niques, our proposed solution aims to reduce the annotation*
 010 *time from hours to mere minutes. However using a direct*
 011 *keypoint detection strategy based on segmentation presents*
 012 *certain issues and does not converge. We test this method,*
 013 *analyze the results and propose methods to overcome the*
 014 *problems.*

1. Introduction

016 Estimating the body conditions of whales holds significant
 017 importance for marine biologists and conservationists, fa-
 018 cilitating assessments of health, movement patterns, and
 019 responses to environmental factors. Traditional methods
 020 like whaling and direct inspection are not only risky but
 021 also costly. However, the emergence of drone technology
 022 enables the non-invasive collection of vast aerial imagery,
 023 providing valuable data for analysis. Currently, annotating
 024 these images manually as shown in Figure 2 consumes con-
 025 siderable time. While existing solutions like Mask-RCNN
 026 [2] offer automation, their dependency on GPU inference
 027 presents practical challenges, particularly in remote field
 028 settings where GPU-enabled devices are impractical. Ad-
 029 dressing this, we propose a novel approach to enable CPU
 030 inference without sacrificing accuracy, leveraging models
 031 optimized for CPU execution. This innovation promises to
 032 streamline annotation processes, empowering marine biolo-
 033 gists with efficient tools for studying whale populations and
 034 aiding in conservation efforts.

2. Related Works

035 The current landscape of whale body mass estimation
 036 requires addressing the need for low-compute solutions
 037 tailored to the constraints of marine biologist field-
 038 work. YOLO[9], renowned for its low-compute efficiency,
 039 emerges as a promising candidate for real-time on-spot de-
 040 tection, as evidenced in recent studies on marine life object
 041 detection[14] [3] [4].

042 Presently, machine learning-based approaches[2] for
 043 body condition estimation typically involve two models:
 044 image segmentation for identifying body contours and key-
 045 point detection for locating the body axis. Combined results
 046 yield width values used by biologists for estimating body
 047 condition parameters. Notably, existing solutions achieve
 048 high accuracy, with segmentation mask accuracy at 96%
 049 and axis detection at 93%. However, training these mod-
 050 els is hindered by the laborious manual annotation process.
 051 For instance, in this study on Southern right whale species
 052 involved manual annotation of 300 images containing 468
 053 whale instances, underscoring the time-consuming nature
 054 of data preparation. While the current inference model
 055 utilizes Detectron2[5], leveraging the Mask-RCNN model,
 056 its dependency on GPU at inference poses practical chal-
 057 lenges in field settings. To mitigate this, transitioning to
 058 YOLO[14] offers a promising alternative.

059 Furthermore, insights from human pose detection re-
 060 search suggest that regression-based [11] approaches offer
 061 faster processing, albeit with slight compromises in accu-
 062 racy and mean Average Precision (mAP) as compared to
 063 one-hot mask[5] and heatmap [8] approaches. Leveraging
 064 these insights, the study explores a detection-localized seg-
 065 mentation method followed by a regression head for key-
 066 point detection. For detection and bounding box drawing,
 067 YOLO[9] serves as a robust choice, while SOTA seg-
 068 mentation methods like Segment Anything (SAM)[6] by Meta
 069 showcase adaptability and superior performance in diverse
 070 domains of computer vision and image processing, promis-
 071 ing advancements in whale body mass estimation tech-
 072 niques. For localized segmentation with a minimal amount
 073 of data, UNet[10] has been proven as a good method. Wu
 074

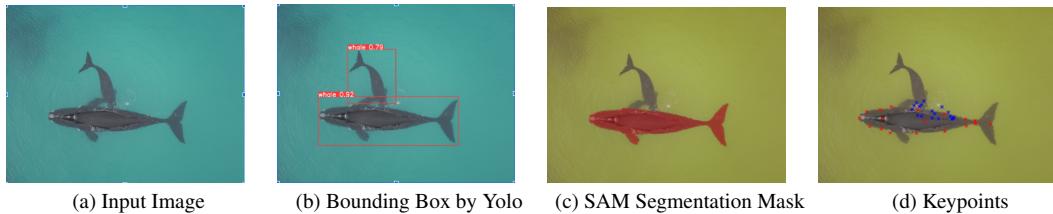


Figure 1. Output at each step. In (d) ground truth is in red and model's output in blue

075 et al.[13] used a UNET architecture to detect and segment
 076 very small objects[10] [7], suggesting its efficiency in seg-
 077 menting small local objects with greater accuracy, including
 078 marine life[12].

079 3. Dataset

080 The dataset comprises approximately 12,000 unlabelled raw
 081 whale images depicting four distinct species. These images
 082 exhibit resolutions of either (2592 x 4608) or (3456
 083 x 4608) pixels. Among these, 400 images featuring South-
 084 ern right whales have been meticulously annotated. A rep-
 085 resentative image from the dataset is presented in Figure
 086 2. In approximately 90% of the images, a pair of whales
 087 is discernible, typically comprising a larger maternal whale
 088 accompanied by a smaller calf. Single whale occurrences
 089 constitute around 4% of the dataset, while the remaining
 090 images portray more than two whales.

091 The ground truth annotations consist of 38 coordinate
 092 values, encompassing (x, y) coordinates for each of the 19
 093 keypoints characterizing the whale's anatomy (refer to Fig-
 094 ure 3). In instances where multiple whales are present in an
 095 image, not all whales are annotated, primarily due to certain
 096 keypoints being obscured from view for some individuals.

097 4. Methodology

098 4.1. Data Preprocessing

099 For the training of the YOLOv8 instance detector tailored
 100 for whale detection, our approach comprised several key
 101 steps. Initially, bounding boxes were manually drawn
 102 around whale instances within 400 images using the CVAT
 103 online annotation tool. Subsequently, these annotations
 104 were converted into YOLO format labels to facilitate train-
 105 ing.

106 Given the absence of keypoint annotations for all whale
 107 instances within an image, we devised a refinement process
 108 to enhance the dataset's quality. This involved the removal
 109 of SAM segmentations for instances lacking keypoint an-
 110 notations, employing a technique known as *closest distance*
 111 *matching*. The following steps outline this process:

- 112 1. Find the center points of all the bounding boxes.
- 113 2. Get the center of the body of the whales from the an-
 114 notations.

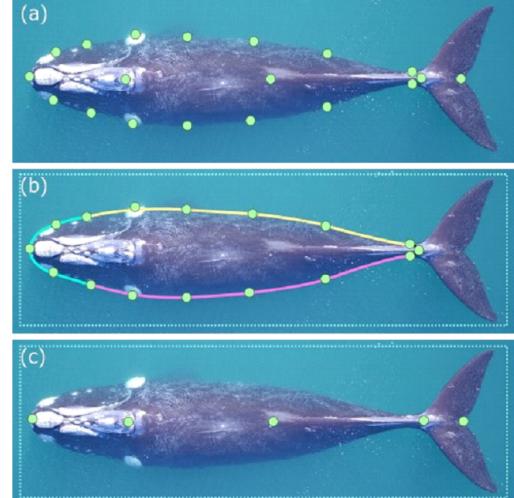


Figure 2. : (a) Point marked by annotators on each whale, (b) 3 cubic splines are fitted on the annotated points to serve as input to the segmentation model that finds the body contour, (c) 5 points along the axis that are used as input to the keypoint detection model to find the body axis. [2]



121 **4.2. Instance Detection**

122 In consideration of real-time testing requirements and the
 123 absence of universal keypoint labels, we bifurcated the
 124 training approach into instance detection and keypoint de-
 125 tector phases. Instance detection, accomplished through
 126 training a YOLOv8 model separately, aimed at detecting
 127 all whale instances within an image during real-time test-
 128 ing. As keypoint labels were unavailable for every whale
 129 instance, image labeling was performed using an online an-
 130 notator, facilitating model training.

131 **4.3. SAM Segmentation**

132 Furthermore, we leveraged the YOLOv8 detector’s bound-
 133 ing box coordinates to fine-tune SAM-generated segmen-
 134 tation masks. This refinement process involved extract-
 135 ing the most confident mask and we excluded masks for the
 136 instances lacking keypoint annotations. Notably, this mask
 137 removal step was confined to the training phase and not in-
 138 tended for real-time application.

139 **4.4. UNet Based Keypoint Detection**

140 A segmentation-based model was used for the keypoint
 141 detection because identifying pixel level features was re-
 142 quired for identifying the keypoints and a segmentaion-
 143 based model could create segmentations localized around
 144 the corresponding keypoints. We created a UNet architec-
 145 ture with four convolution layers of filters (32,64,128,256),
 146 a strided convolution layer with 512 filters and then corre-
 147 sponding de-convolution layers in the decoder step. A sig-
 148 moid layer was added after the decoder step. This UNet
 149 backbone produced an output of the same 2D dimensions
 150 as that of the input image. To obtain keypoint regressions
 151 from the this output, we added an MLP layer of layer sizes
 152 (1024,256,38) to output 38 keypoint values.

153 For the input of the UNet model, we used the raw whale im-
 154 age and concatenated the segmentation mask as the fourth
 155 channel, thus created an input of dimensions $(4, H, W)$.
 156 The segmentation mask was added to provide context and
 157 the boundary of the whale as an additional information to
 158 help in localization. So if an image had two annotations
 159 for two distinct labels, then by adding masks for different
 160 whales separately to that image, we created two different
 161 training inputs, one for each annotated whale. All images
 162 were of very high resolution and did not have the same di-
 163 mensions, thus making it difficult to load them into memory
 164 and pass as an input to the model. Thus we zero-padded the
 165 smaller images to match the sizes to (3456, 4608) and then
 166 resized them using bilinear downsampling to one-eighth of
 167 their original size (432, 576).

168 The overall architecture is shown in Figure 4.

169 **4.5. Loss**

170 For the loss, we used **Mean Squared Error** measured be-
 171 tween the predicted coordinate values and the actual coor-
 172 dinate values obtained from the annotations.

173 **5. Experiments**174 **5.1. YOLO Detection**

175 We initially evaluated instance detection using the YOLOv8
 176 model with pretrained weights. However, the model fre-
 177 quently failed to identify the tail fin of the whale as a part
 178 of the body and occasionally labeled the tail as a separate
 179 instance. To address this, we manually created box labels
 180 and custom-trained the YOLOv8 model using 350 exam-
 181 ples, allocating 300 for training and 50 for validation. Train-
 182 ing occurred over 100 epochs with default learning param-
 183 eters. YOLOv8 automatically logged evaluation data for
 184 both training and validation examples.

185 Subsequently, we conducted qualitative testing on 50 ad-
 186 ditional test examples. Through fine-tuning, the model ex-
 187 hibited significant improvement, producing tight and pre-
 188 cise bounding boxes around whales with high confidence
 189 levels.

190 **5.2. SAM Segmentation**

191 We did not train the SAM model explicitly on our data since
 192 we did not have segmentation annotations. However, we as-
 193 sisted the SAM model with the bounding box coordinates
 194 received from the previous stage to localize the segmenta-
 195 tions. With the box prompts, the model could generate pre-
 196 cise masks which we evaluated with the qualitative analysis
 197 of random samples. The segmentations of the downsampled
 198 images were also manually observed to check for any dis-
 199 tortions. The masks were still significantly precise, however
 200 there were color distortions.

201 **5.3. UNet based Keypoint Detection**202 **Training**

203 The SAM generated masks were added as the fourth chan-
 204 nel to the images and then downsampled. We divided the
 205 data into train, validation and test sets in the split ratio of
 206 76:9:15 respectively.

207 Even with the downsampled images, we could not store
 208 the entire training set into the memory at a time. Hence we
 209 divided the training set into six sets and successively trained
 210 them for 70 epochs each. After training one set, we saved
 211 the parameters and used them for the next set, thus not train-
 212 ing every set from scratch. Initially we used an AdamW
 213 optimizer, with a leraning rate of 0.001 and weight decay
 214 of 0.00001. In the initial stages of training, the training
 215 loss and testing loss decreased gradually. However eventu-
 216 ally the drop of validation loss started reducing, indicating

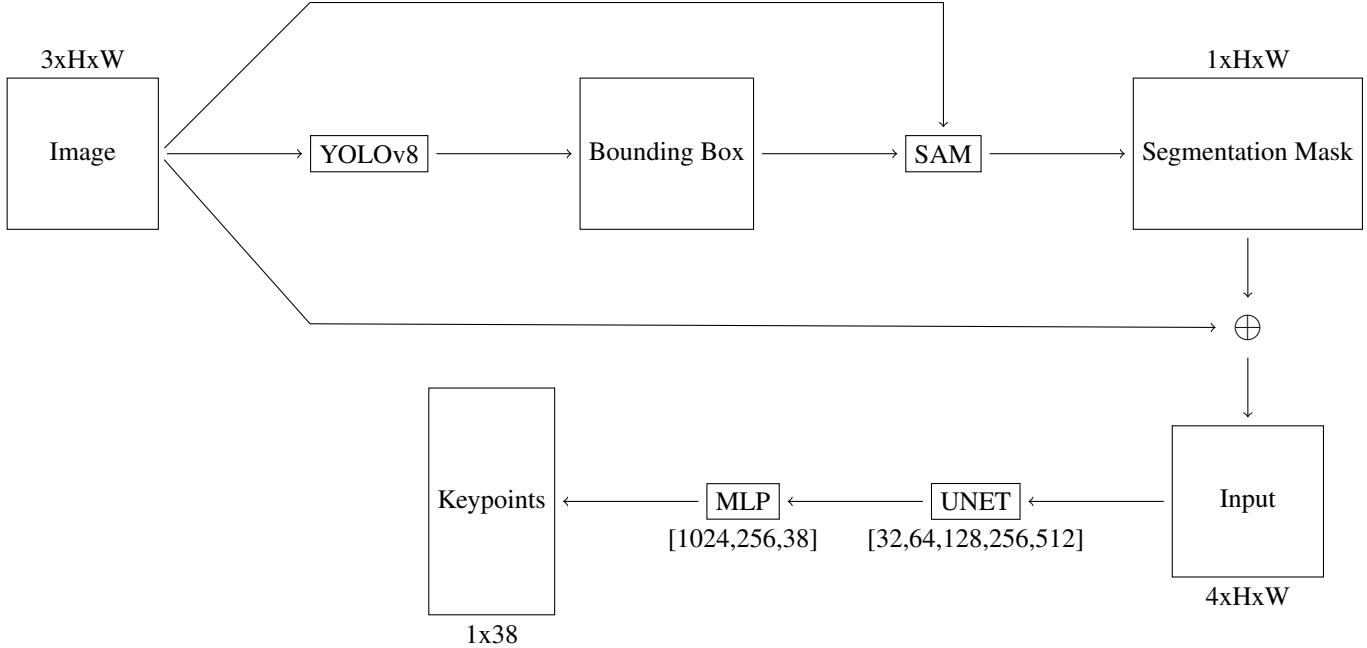


Figure 4. Architecture

217 overfitting (refer Figure 5). Thus in later stages, we exper-
 218 imented with decreasing the learning rates upto 0.000005
 219 and increasing the weight decay upto 0.001. However there
 220 was no visible improvement and the validation loss stag-
 221 nated in all scenarios.

222 We also tested the training by reducing the downsam-
 223 pling ratio from 8 to 4. There were no noticeable changes
 224 in the learning curve pattern, thus discarding the possibility
 225 of downsampling causing the poor predictions.
 226

227 Keypoints Evaluation

228 On an average, each keypoint was about 100 pixels off from
 229 the actual values. We also visually compared the positions
 230 of predicted keypoints and actual locations. Even though
 231 the predicted keypoints were aggregated over the whale,
 232 which probably happened because of the segmentation
 233 mask provided as input, there was no recognizable pattern.
 234 This indicated towards any lack of feature level under-
 235 standing of the keypoints. In some images, the shapes formed
 236 by the keypoints resembled the shape of a whale, however
 237 such instances were rare and resulted from overfitting.
 238

239 UNet Block Outputs

240 We also examined the outputs generated by the UNet block
 241 before the MLP layer. Using the MSE loss with respect
 242 to the keypoints, we expected the concentration of distinct
 243 segmentations around the keypoints. Contrary to our
 244 expectations, the generated masks consisted segmentations
 245 over the entire body of the whale (refer Figure 8 Column

246 4). These masks were similar to the ones created by the
 247 SAM model which we had concatenated in the input. Thus,
 248 the intended effect of UNet was not observed and we could
 249 have obtained similar results by directly using SAM masks
 250 with the MLP layer.

251 YOLOv8 Pretrained Detection

252 We tested the keypoint detection using off-the-shelf
 253 YOLOv8 keypoint detector. We trained the model for 75
 254 epochs with the same train and test split that we used for
 255 YOLOv8 instance detection. Though we saw incredible re-
 256 sults previously when training the bounding box prediction
 257 individually, the results for the keypoint detection were bad
 258 (refer Figure 9) . Even the bounding box results were not
 259 good and whale instances were missed in many examples.
 260 This happened because of the lack of consistencies in the
 261 annotations available for all whales present in an image. As
 262 a result, the whale detection deteriorated.

264 6. Conclusion

265 Due to the limitations with available annotations, we de-
 266 veloped a model for keypoint detection with separate in-
 267 stance detector, segmentation mask generator and keypoint
 268 detector. We used a regression head on the top of the
 269 segmentation backbone to create localized patches around the
 270 keypoints and then intended to use them to infer the coordi-
 271 nates. However the UNet created full whale segmentations,
 272 similar to ones created by the SAM, thus deeming the addi-
 273

Submission .

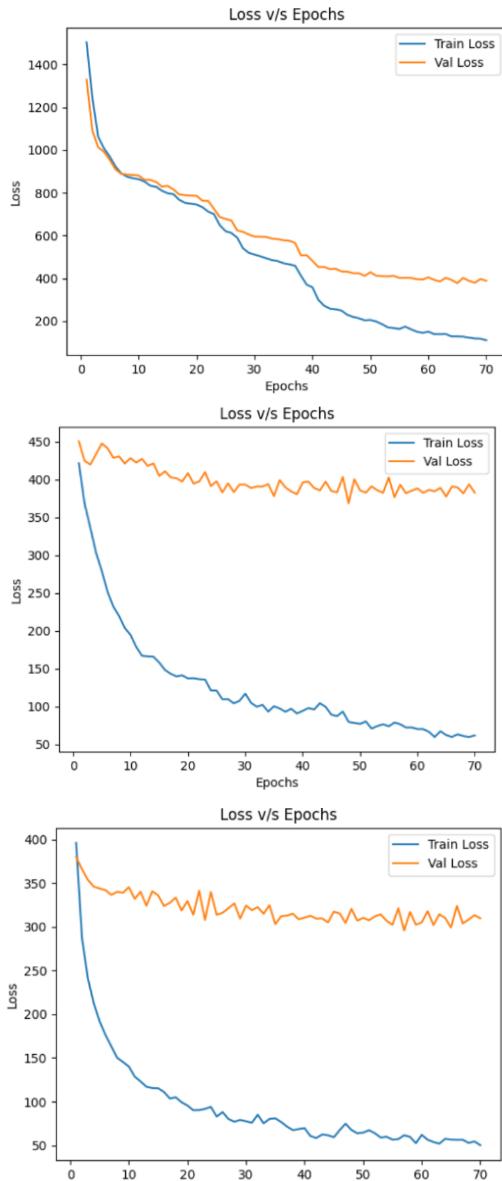


Figure 5. Training Loss and Validation Loss v/s Epochs for UNet based Keypoint Detection

274

tion of the UNet block redundant.

275

7. Future Works

276

For creating a top-down keypoint detection model based on the segmentation, we need to develop a model which is explicitly trained to generate localized patches in form of heatmaps, instead of training them to predict very specific pixel coordinates. Segmentation model cannot converge to single pixel level prediction. In such methods, we generate a heatmap which is basically a gaussian mask with very low standard deviation and centered at the keypoints, so

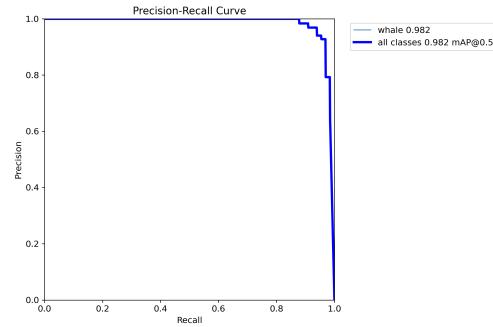


Figure 6. PR Curve for YOLOv8 Instance Detection

it becomes a small patch. Then we train the segmentation model to predict those patches, instead of a single pixel keypoint. Some non-maximal suppression strategy can be used to infer keypoints from those concentrated patches. We will also explore methods for network training refinement and post processing [1] to improve the predictions. Data augmentations can also be used to increase the dataset size and variations.

284
285
286
287
288
289
290
291
292
293
294
295
296

Finally, we will build a human-in-the-loop model to rectify the outputs using human supervision. We will determine ways to perform active learning sampling using uncertainty and diversity sampling strategies.

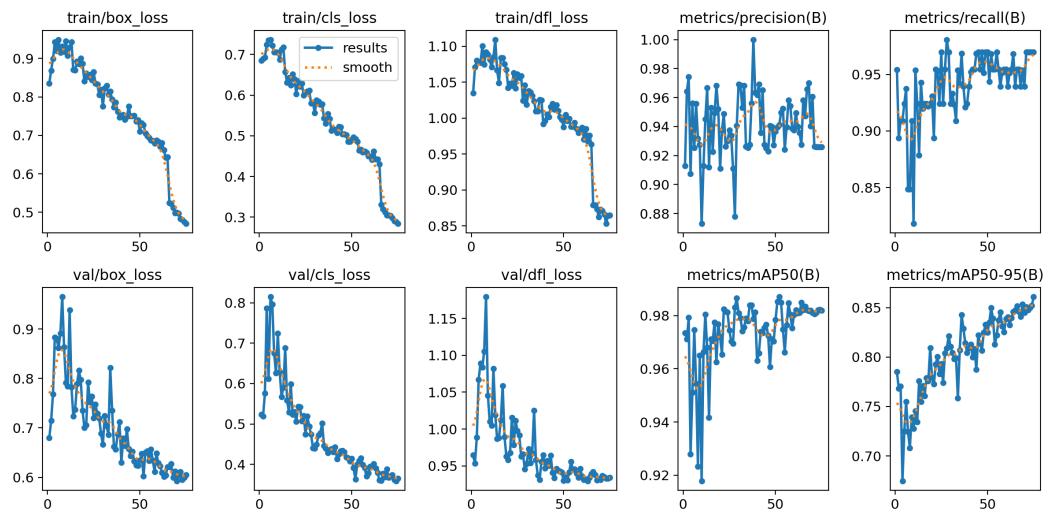


Figure 7. Train and Validation loss v/s Epochs and Mean Average Precision Scores for YOLOv8 Instance Detection



Figure 8. Keypoint Detection Results. (Order: Input-> YOLO-> SAM-> UNET-> MLP)



Figure 9. Keypoints Detection with YOLOv8 Keypoints Detector

297 **References**

- 298 [1] Haoming Chen, Runyang Feng, Sifan Wu, Hao Xu,
299 Fengcheng Zhou, and Zhenguang Liu. 2d human pose es-
300 timation: A survey, 2022. 5
- 301 [2] Duncan Irschick Fredrik Christiansen Chhandak Bagchi,
302 Josh Medina. An automated tool to extract body morpho-
303 metric data from drone aerial photographs of southern right
304 whales to measure body condition. Technical report, Univer-
305 sity of Massachusetts and Aarhus University, 2023. 1, 2
- 306 [3] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian
307 Sun. YOLOX: exceeding YOLO series in 2021. *CoRR*,
308 abs/2107.08430, 2021. 1
- 309 [4] An Guo, Kaiqiong Sun, and Ziyi Zhang. A lightweight
310 yolov8 integrating fasternet for real-time underwater object
311 detection. *Journal of Real-Time Image Processing*, 21, 2024.
312 1
- 313 [5] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Gir-
314 shick. Mask r-cnn, 2018. 1
- 315 [6] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao,
316 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer White-
317 head, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and
318 Ross Girshick. Segment anything, 2023. 1
- 319 [7] Martin Kolarík, Radim Burget, Václav Uher, Kamil Říha,
320 and Malay Kishore Dutta. Optimized high resolution 3d
321 dense-u-net network for brain and spine segmentation. *Ap-
322 plied Sciences*, 2019. 2
- 323 [8] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked
324 hourglass networks for human pose estimation. *CoRR*,
325 abs/1603.06937, 2016. 1
- 326 [9] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick,
327 and Ali Farhadi. You only look once: Unified, real-time ob-
328 ject detection. *CoRR*, abs/1506.02640, 2015. 1
- 329 [10] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net:
330 Convolutional networks for biomedical image segmentation.
331 *CoRR*, abs/1505.04597, 2015. 1, 2
- 332 [11] Xiao Sun, Bin Xiao, Shuang Liang, and Yichen Wei. Integral
333 human pose regression. *CoRR*, abs/1711.08229, 2017. 1
- 334 [12] Lidiya Thampi, Rosemol Thomas, Suraj Kamal, Arun A.
335 Balakrishnan, T. P. Mithun Haridas, and M. H. Supriya.
336 Analysis of u-net based image segmentation model on under-
337 water images of different species of fishes. In *2021 Interna-
338 tional Symposium on Ocean Technology (SYMPOL)*, pages
339 1–5, 2021. 2
- 340 [13] Xin Wu, Danfeng Hong, and Jocelyn Chanussot. Uiu-net: U-
341 net in u-net for infrared small object detection. *IEEE Trans-
342 actions on Image Processing*, 32:364–376, 2023. 2
- 343 [14] J. Zhong, M. Li, J. Qin, Y. Cui, K. Yang, and H. Zhang. Real-
344 time marine animal detection using yolo-based deep learn-
345 ing networks in the coral reef ecosystem. *The International
346 Archives of the Photogrammetry, Remote Sensing and Spa-
347 tial Information Sciences*, XLVI-3/W1-2022:301–306, 2022.
348 1