# Detecting Marine Species in Echograms via Traditional, Hybrid, and Deep Learning Frameworks

Tunai Porto Marques \*1, Alireza Rezvanifar \*1, Melissa Cote 1, Alexandra Branzan Albu 1, Kaan Ersahin 2, Todd Mudge 2, Stephane Gauthier 3 \*: equal contribution

> <sup>1</sup> University of Victoria, BC, Canada <sup>2</sup> ASL Environmental Sciences, Victoria, Canada <sup>3</sup> Institute of Ocean Sciences, Fisheries and Oceans Canada, Sidney, BC, Canada



25th International Conference on Pattern Recognition. Milan, Italy, January 10-15 2021

1 Motivation

(4)

- Non-invasive
- Thorough

- Non-invasive
- Thorough

Echosounders: ASL Environmental Science's AZFP [1]





- Non-invasive
- Thorough

Echosounders: ASL Environmental Science's AZFP [1]

Manual or semi-automatic methods:

- Time-consuming
- Prone to inter-expert disagreements



Illustration of an AZFP-generated echogram [1].



- Non-invasive
- Thorough

Echosounders: ASL Environmental Science's AZFP [1]

Manual or semi-automatic methods:

- Time-consuming
- Prone to inter-expert disagreements







- Non-invasive
- Thorough

Echosounders: ASL Environmental Science's AZFP [1]

Manual or semi-automatic methods:

- Time-consuming
- Prone to inter-expert disagreements





Perform a comparison between machine learning-based approaches for the automatic detection of marine species in echograms.

Motivation

- Non-invasive
- Thorough

Echosounders: ASL Environmental Science's AZFP [1]

Manual or semi-automatic methods:

- Time-consuming
- Prone to inter-expert disagreements





Perform a comparison between machine learning-based approaches for the automatic detection of marine species in echograms.

Motivation

Multi-frequency, heterogeneous morphology



Sample 1-hour echograms measured with a fixed-position echosounder (see Sec. III-A). First row: (a) 67 kHz echogram with range (y-) and time (x-) axes, and acoustic echo intensity; (b)-(d) the same echogram at other frequencies. Second row: four echograms (only 67 kHz is shown) containing various schools identified by red bounding boxes.

1 Motivation

| Approach        | Localization         | Classification      |
|-----------------|----------------------|---------------------|
| 1) Hand-crafted | Custom ROI extractor | ML image classifier |
| 2) Hybrid       | Custom ROI extractor | DL image classifier |
| 3) End-to-end   | DL object detector   | DL object detector  |

ROI = region of interest, ML = machine learning, DL = deep learning

2 Proposed Approach Results

(4)

#### Hand-crafted

- Most common in the related literature (see Sec. I-B from the manuscript)
- Novel Region of Interest (ROI) extractor
  - Herring-specific assumptions
- Support Vector Machine (SVM)
  Four hand-crafted features

| Approach        | Localization         | Classification      |
|-----------------|----------------------|---------------------|
| 1) Hand-crafted | Custom ROI extractor | ML image classifier |
| 2) Hybrid       | Custom ROI extractor | DL image classifier |
| 3) End-to-end   | DL object detector   | DL object detector  |

ROI = region of interest, ML = machine learning, DL = deep learning

#### Hand-crafted

- Most common in the related literature (see Sec. I-B from the manuscript)
- Novel Region of Interest (ROI) extractor
  - Herring-specific assumptions
- Support Vector Machine (SVM)
  Four hand-crafted features

#### Hybrid

- Novel Region of Interest (ROI) extractor
  - Herring-specific assumptions
- Custom deep learning-based image classifiers (ResNet [22], DenseNet[23], Inception [24])

| Approach        | Localization         | Classification      |
|-----------------|----------------------|---------------------|
| 1) Hand-crafted | Custom ROI extractor | ML image classifier |
| 2) Hybrid       | Custom ROI extractor | DL image classifier |
| 3) End-to-end   | DL object detector   | DL object detector  |

ROI = region of interest, ML = machine learning, DL = deep learning

#### Hand-crafted

- Most common in the related literature (see Sec. I-B from the manuscript)
- Novel Region of Interest (ROI) extractor
  - Herring-specific assumptions
- Support Vector Machine (SVM)
  Four hand-crafted features

#### Hybrid

- Novel Region of Interest (ROI) extractor
  - Herring-specific assumptions
- Custom deep learning-based image classifiers (ResNet [22], DenseNet[23], Inception [24])

| Approach        | Localization         | Classification      |
|-----------------|----------------------|---------------------|
| 1) Hand-crafted | Custom ROI extractor | ML image classifier |
| 2) Hybrid       | Custom ROI extractor | DL image classifier |
| 3) End-to-end   | DL object detector   | DL object detector  |

ROI = region of interest, ML = machine learning, DL = deep learning

#### End-to-end

- Custom-trained end-to-end object detectors (Faster R-CNN [25], YOLOv2 [26])
- Novel tiling approach

 $\bigcirc$ 

Proposed Approach Results



(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.



(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.

"Counts" Images (4 frequency channels)





Echograms from the same site captured using four different frequencies.

on Proposed Approach Results Di

(4)

"Counts" Images (4 frequency channels)



- Median filter on the x-axis (time direction)
- Short signals (individual fish)





• Region-specific thresholding





• Removes small sets of connected components





"Counts" Images (4 frequency channels)



• Only positions with intensities > 2 are kept (i.e. signals present in three or more frequencies)

3.5

2.5

1.5

0.5



2 Proposed Approach

Filtered Score Matrix (left), unfiltered Score Matrix (right).

"Counts" Images (4 frequency channels) 67 kHz 125 kHz kHz 455 kHz kHz kHz Denoising 67 kHz 125 125 200 200 455 kHz kHz kHz kHz kHz Adaptive Thresholding 67 kHz 125 kHz kHz 455 kHz 455 kHz **Opening + Closing** Score Matrix Filtering by size Filtering by orientation Detected **ROIs** 

• Only positions with intensities > 2 are kept (i.e. signals present in three or more frequencies)



"Counts" Images (4 frequency channels)



• Removes small acoustic responses (usually associated with individual fish)

2 Proposed Approach

• Threshold: 50 pixels



"Counts" Images (4 frequency channels)



- Bounding boxes (BBs) around remaining structures.
- 67 kHz (more pronounced response for schools of herring)



"Counts" Images (4 frequency channels)



- Ignore a BB with  $\alpha < 60^{\circ}$
- Contextual information about schools of herring



Bonding

2 Proposed Approach α

Output of BB orientation filtering (final ROIs).







(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.



(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.

### Hand-crafted (SVM) approach

#### Output from the ROI extractor

![](_page_27_Figure_2.jpeg)

1 2 3 otivation Proposed Approach Results Discu

### Hand-crafted (SVM) approach

![](_page_28_Picture_1.jpeg)

1 (2) (3) (4) otivation Proposed Approach Results Discussion

![](_page_29_Figure_0.jpeg)

(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.

![](_page_30_Figure_0.jpeg)

(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.

### Hybrid approach

![](_page_31_Figure_1.jpeg)

![](_page_32_Figure_0.jpeg)

(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.

![](_page_33_Figure_0.jpeg)

(b)

General block diagram of hand-crafted and hybrid approaches (a) and end-to-end approach (b). Red bounding boxes highlight the identified targets in the final results, shown on colour-coded echograms for better visualization.

![](_page_34_Picture_1.jpeg)

![](_page_34_Picture_2.jpeg)

Scale **a** 

![](_page_35_Picture_1.jpeg)

![](_page_35_Picture_2.jpeg)

![](_page_36_Figure_1.jpeg)

A generic convolution-based feature extractor (adapted from Pinaya et al., 2020).

Pinaya, W. H. L. et al. Machine Learning: Methods and Applications to Brain Disorders. Academic Press, 2020, pp. 173-191.

![](_page_37_Figure_1.jpeg)

A generic convolution-based feature extractor (adapted from Pinaya et al., 2020).

Pinaya, W. H. L. et al. Machine Learning: Methods and Applications to Brain Disorders. Academic Press, 2020, pp. 173-191.

(2)

Proposed Approach Results

## Tiling strategy

![](_page_38_Picture_1.jpeg)

![](_page_38_Picture_2.jpeg)

### Tiling strategy

![](_page_39_Figure_1.jpeg)

Pinaya, W. H. L. et al. Machine Learning: Methods and Applications to Brain Disorders. Academic Press, 2020, pp. 173-191.

(2)

Proposed Approach

(4)

### Tiling strategy

![](_page_40_Figure_1.jpeg)

# (a) Original echogram

# (b) Three overlapping tiles

![](_page_40_Figure_4.jpeg)

# (c) Multiple detection (inference)

![](_page_40_Picture_6.jpeg)

# (d) Concatenated detection results

(2) Proposed Approach

Tiling strategy used for training and inference. Original echogram (571x1200) with ground truth annotation highlighted in yellow (a). Tiles of 340x340 created around the annotation for the training phase, highlighted in yellow (b). Multiple tiles (green), with the first tile highlighted in black, and individual detection results (red bounding boxes) obtained during inference (c). Detection bounding boxes (red) obtained after post-processing (d).

### Dataset

![](_page_41_Figure_1.jpeg)

- Canada's Department of Fisheries and Oceans (DFO)
- ASL Environmental Sciences' AZFP [1] echosounder
- Okisollo channel, BC, Canada
- May October 2015 and 2016

![](_page_41_Picture_6.jpeg)

### Dataset

![](_page_42_Figure_1.jpeg)

vation Proposed

Results

![](_page_43_Figure_0.jpeg)

## Qualitative results

#### 1) Hand-crafted (SVM)

![](_page_43_Picture_3.jpeg)

#### 2) Hybrid (InceptionV3 [24])

3) End-to-end (YOLOv2 [26])

![](_page_43_Figure_6.jpeg)

![](_page_43_Picture_7.jpeg)

## Quantitative results

| Method                            | IoU | Р     | R     | F1    |
|-----------------------------------|-----|-------|-------|-------|
| 1) Hand-crafted (SVM [29])        | 0.3 | 45.41 | 65.63 | 53.67 |
| 1) Hand-crafted (SVM [29])        | 0.5 | 41.62 | 60.16 | 49.20 |
| 2) Hybrid (ResNet-50 [22])        | 0.3 | 72.37 | 85.94 | 78.57 |
| 2) Hybrid (ResNet-50 [22])        | 0.5 | 62.50 | 74.20 | 67.80 |
| 2) Hybrid (DenseNet-201 [23])     | 0.3 | 64.37 | 87.50 | 74.17 |
| 2) Hybrid (DenseNet-201 [23])     | 0.5 | 55.75 | 75.78 | 64.24 |
| 2) Hybrid (InceptionV3 [24])      | 0.3 | 78.99 | 85.16 | 81.95 |
| 2) Hybrid (InceptionV3 [24])      | 0.5 | 68.12 | 73.44 | 70.68 |
| 3) End-to-end (YOLOv2 [26])       | 0.3 | 73.08 | 89.06 | 80.28 |
| 3) End-to-end (YOLOv2 [26])       | 0.5 | 67.31 | 82.03 | 73.94 |
| 3) End-to-end (Faster R-CNN [25]) | 0.3 | 66.90 | 74.22 | 70.37 |
| 3) End-to-end (Faster R-CNN [25]) | 0.5 | 45.07 | 50    | 47.41 |

P = precision, R = recall, F1 = F1-score

3

Motivation Proposed Approach Results

(4)

## Quantitative results

| Method                            | IoU | P     | R     | <b>F1</b> | - |
|-----------------------------------|-----|-------|-------|-----------|---|
| 1) Hand-crafted (SVM [29])        | 0.3 | 45.41 | 65.63 | 53.67     | - |
| 1) Hand-crafted (SVM [29])        | 0.5 | 41.62 | 60.16 | 49.20     | - |
| 2) Hybrid (ResNet-50 [22])        | 0.3 | 72.37 | 85.94 | 78.57     | - |
| 2) Hybrid (ResNet-50 [22])        | 0.5 | 62.50 | 74.20 | 67.80     | - |
| 2) Hybrid (DenseNet-201 [23])     | 0.3 | 64.37 | 87.50 | 74.17     | - |
| 2) Hybrid (DenseNet-201 [23])     | 0.5 | 55.75 | 75.78 | 64.24     | - |
| 2) Hybrid (InceptionV3 [24])      | 0.3 | 78.99 | 85.16 | 81.95     |   |
| 2) Hybrid (InceptionV3 [24])      | 0.5 | 68.12 | 73.44 | 70.68     | - |
| 3) End-to-end (YOLOv2 [26])       | 0.3 | 73.08 | 89.06 | 80.28     |   |
| 3) End-to-end (YOLOv2 [26])       | 0.5 | 67.31 | 82.03 | 73.94     | - |
| 3) End-to-end (Faster R-CNN [25]) | 0.3 | 66.90 | 74.22 | 70.37     | - |
| 3) End-to-end (Faster R-CNN [25]) | 0.5 | 45.07 | 50    | 47.41     | _ |

P = precision, R = recall, F1 = F1-score

3

Motivation Proposed Approach Results

## Quantitative results

| Method                            | IoU | P     | R     | <b>F1</b> |  |
|-----------------------------------|-----|-------|-------|-----------|--|
| 1) Hand-crafted (SVM [29])        | 0.3 | 45.41 | 65.63 | 53.67     |  |
| 1) Hand-crafted (SVM [29])        | 0.5 | 41.62 | 60.16 | 49.20     |  |
| 2) Hybrid (ResNet-50 [22])        | 0.3 | 72.37 | 85.94 | 78.57     |  |
| 2) Hybrid (ResNet-50 [22])        | 0.5 | 62.50 | 74.20 | 67.80     |  |
| 2) Hybrid (DenseNet-201 [23])     | 0.3 | 64.37 | 87.50 | 74.17     |  |
| 2) Hybrid (DenseNet-201 [23])     | 0.5 | 55.75 | 75.78 | 64.24     |  |
| 2) Hybrid (InceptionV3 [24])      | 0.3 | 78.99 | 85.16 | 81.95     |  |
| 2) Hybrid (InceptionV3 [24])      | 0.5 | 68.12 | 73.44 | 70.68     |  |
| 3) End-to-end (YOLOv2 [26])       | 0.3 | 73.08 | 89.06 | 80.28     |  |
| 3) End-to-end (YOLOv2 [26])       | 0.5 | 67.31 | 82.03 | 73.94 <   |  |
| 3) End-to-end (Faster R-CNN [25]) | 0.3 | 66.90 | 74.22 | 70.37     |  |
| 3) End-to-end (Faster R-CNN [25]) | 0.5 | 45.07 | 50    | 47.41     |  |

P = precision, R = recall, F1 = F1-score

3

Motivation Proposed Approach Results

4

## Discussion

- Three machine learning-based approaches:
  - Hand-crafted features
  - Hybrid
  - End-to-end
- ROI extractor
- Tiling strategy (scale limitations)
- Dataset of 358 annotated echograms
- Hybrid approach (IoU threshold 0.3)
- End-to-end approach (IoU threshold 0.5)

![](_page_47_Picture_10.jpeg)

![](_page_47_Figure_11.jpeg)

![](_page_47_Picture_12.jpeg)

## Discussion

- Three machine learning-based approaches:
  - Hand-crafted features
  - Hybrid
  - End-to-end
- ROI extractor
- Tiling strategy (scale limitations)
- Dataset of 358 annotated echograms
- Hybrid approach (IoU threshold 0.3)
- End-to-end approach (IoU threshold 0.5)

![](_page_48_Picture_10.jpeg)

![](_page_48_Picture_11.jpeg)

![](_page_48_Figure_12.jpeg)

![](_page_48_Picture_13.jpeg)

# Thank you!

![](_page_49_Picture_1.jpeg)

# References

- D. Lemon, P. Johnston, J. Buermans, E. Loos, G. Borstad, and L. Brown, "Multiple-frequency moored sonar for continuous observations of zooplankton and fish," in 2012 Oceans. IEEE, 2012, pp. 1–6.
- [2] D. G. Reid, "Report on echo trace classification," ICES Cooperative Research Report, no. 238, 2000.
- [3] T. K. Stanton, "30 years of advances in active bioacoustics: A personal perspective," *Method. Oceanogr.*, vol. 1-2, pp. 49–77, 2012.
- [4] J. K. Horne, "Acoustic approaches to remote species identification: A review," *Fish. Oceanogr.*, vol. 9, no. 4, pp. 356–71, 2000.
- [5] D. Reid, C. Scalabrin, P. Petitgas, J. Masse, R. Aukland, P. Carrera *et al.*, "Standard protocols for the analysis of school based data from echo sounder surveys," *Fish. Res.*, vol. 47, no. 2-3, pp. 125–36, 2000.
- [6] P. LeFeuvre, G. Rose, R. Gosine, R. Hale, W. Pearson, and R. Khan, "Acoustic species identification in the Northwest Atlantic using digital image processing," *Fish. Res.*, vol. 47, no. 2-3, pp. 137–47, 2000.
- [7] A. Charef, S. Ohshimo, I. Aoki, and N. Al Absi, "Classification of fish schools based on evaluation of acoustic descriptor characteristics," *Fish. Sci.*, vol. 76, no. 1, pp. 1–11, 2010.
- [8] A. G. Cabreira, M. Tripode, and A. Madirolas, "Artificial neural networks for fish-species identification," *ICES J. Mar. Sci.*, vol. 66, no. 6, pp. 1119– 29, 2009.
- [9] H. Robotham, P. Bosch, J. C. Gutiérrez-Estrada, J. Castillo, and I. Pulido-Calvo, "Acoustic identification of small pelagic fish species in Chile using support vector machines and neural networks," *Fish. Res.*, vol. 102, no. 1-2, pp. 115–22, 2010.
- [10] S. Gauthier, J. Oeffner, and R. L. O'Driscoll, "Species composition and acoustic signatures of mesopelagic organisms in a subtropical convergence zone, the New Zealand Chatham Rise," *Mar. Ecol. Prog. Ser.*, vol. 503, pp. 23–40, 2014.

- [11] N. G. Fallon, S. Fielding, and P. G. Fernandes, "Classification of Southern Ocean krill and icefish echoes using random forests," *ICES J. Mar. Sci.*, vol. 73, no. 8, pp. 1998–2008, 2016.
- [12] K. Malde, N. O. Handegard, L. Eikvil, and A.-B. Salberg, "Machine intelligence and the data-driven future of marine science," *ICES J. Mar. Sci.*, p. fsz057, 2019.
- [13] Y. Hirama, S. Yokoyama, T. Yamashita, H. Kawamura, K. Suzuki, and M. Wada, "Discriminating fish species by an Echo sounder in a set-net using a CNN," in 21st Asia Pac. Symp. Intell. Evol. Syst. (IES). IEEE, 2017, pp. 112–5.
- [14] Y. Shang and J. Li, "Study on echo features and classification methods of fish species," in 10th Int. Conf. Wirel. Commun. Signal Process. (WCSP). IEEE, 2018, pp. 1–6.
- [15] O. Brautaset, A. U. Waldeland, E. Johnsen, K. Malde, L. Eikvil, A.-B. Salberg *et al.*, "Acoustic classification in multifrequency echosounder data using deep convolutional neural networks," *ICES J. Mar. Sci.*, 2020.
- [16] R. J. Korneliussen, Y. Heggelund, G. J. Macaulay, D. Patel, E. Johnsen, and I. K. Eliassen, "Acoustic identification of marine species using a feature library," *Method. Oceanogr.*, vol. 17, pp. 187–205, 2016.
- [17] L. Liu, H. Lu, Z. Cao, and Y. Xiao, "Counting fish in sonar images," in 25th IEEE Int. Conf. Image Process. (ICIP). IEEE, 2018, pp. 3189–93.
- [18] D. Glukhov, R. Bohush, J. Mäkiö, and T. Hlukhava, "A joint application of fuzzy logic approximation and a deep learning neural network to build fish concentration maps based on sonar data," in 2nd Int. Workshop Comput. Model. Intell. Syst. (CMIS), 2019, pp. 133–42.
- [19] G. French, M. Mackiewicz, M. Fisher, M. Challiss, P. Knight, B. Robinson *et al.*, "JellyMonitor: automated detection of jellyfish in sonar images using neural networks," in *14th IEEE Int. Conf. Signal Process. (ICSP)*. IEEE, 2018, pp. 406–12.

tivation Proposed Approach Resu

Discussion

# References

- [20] H. Wang, Y. Yu, Y. Cai, X. Chen, L. Chen, and Q. Liu, "A comparative study of state-of-the-art deep learning algorithms for vehicle detection," *IEEE Intell. Transp. Syst. Mag.*, vol. 11, no. 2, pp. 82–95, 2019.
- [21] S. H. Kassani and P. H. Kassani, "A comparative study of deep learning architectures on melanoma detection," *Tissue Cell*, vol. 58, pp. 76–83, 2019.
- [22] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–8.
- [23] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 4700–8.
- [24] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *IEEE Conf. Comput.*

Vis. Pattern Recognit. (CVPR), 2016, pp. 2818–26.

- [25] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards realtime object detection with region proposal networks," in Adv. Neural Inf. Process. Syst. (NIPS), 2015, pp. 91–9.
- [26] J. Redmon and A. Farhadi, "YOLO9000: better, faster, stronger," in IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2017, pp. 7263–71.
- [27] A. Rezvanifar, T. P. Marques, M. Cote, A. B. Albu, A. Slonimer, T. Tolhurst *et al.*, "A deep learning-based framework for the detection of schools of herring in echograms," *arXiv preprint arXiv:1910.08215*, 2019.
- [28] D. Bradley and G. Roth, "Adaptive thresholding using the integral image," J. Graph. Tools, vol. 12, no. 2, pp. 13–21, 2007.

- [29] C. Cortes and V. Vapnik, "Support-vector networks," Mach. Learn., vol. 20, no. 3, pp. 273–97, 1995.
- [30] R. Girshick, "Fast R-CNN," in IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2015, pp. 1440–8.
- [31] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2014.
- [32] J. Kiefer and J. Wolfowitz, "Stochastic estimation of the maximum of a regression function," Ann. Math. Statist., vol. 23, no. 3, pp. 462–6, 1952.

Discussion