

Deep learning detection and classification of baleen whale vocalizations using a novel data representation

Mark Thomas^{1,2,*}, Bruce Martin², Katie Kowarski², Briand Gaudet², and Stan Matwin¹

¹Dalhousie University, Faculty of Computer Science ²JASCO Applied Sciences

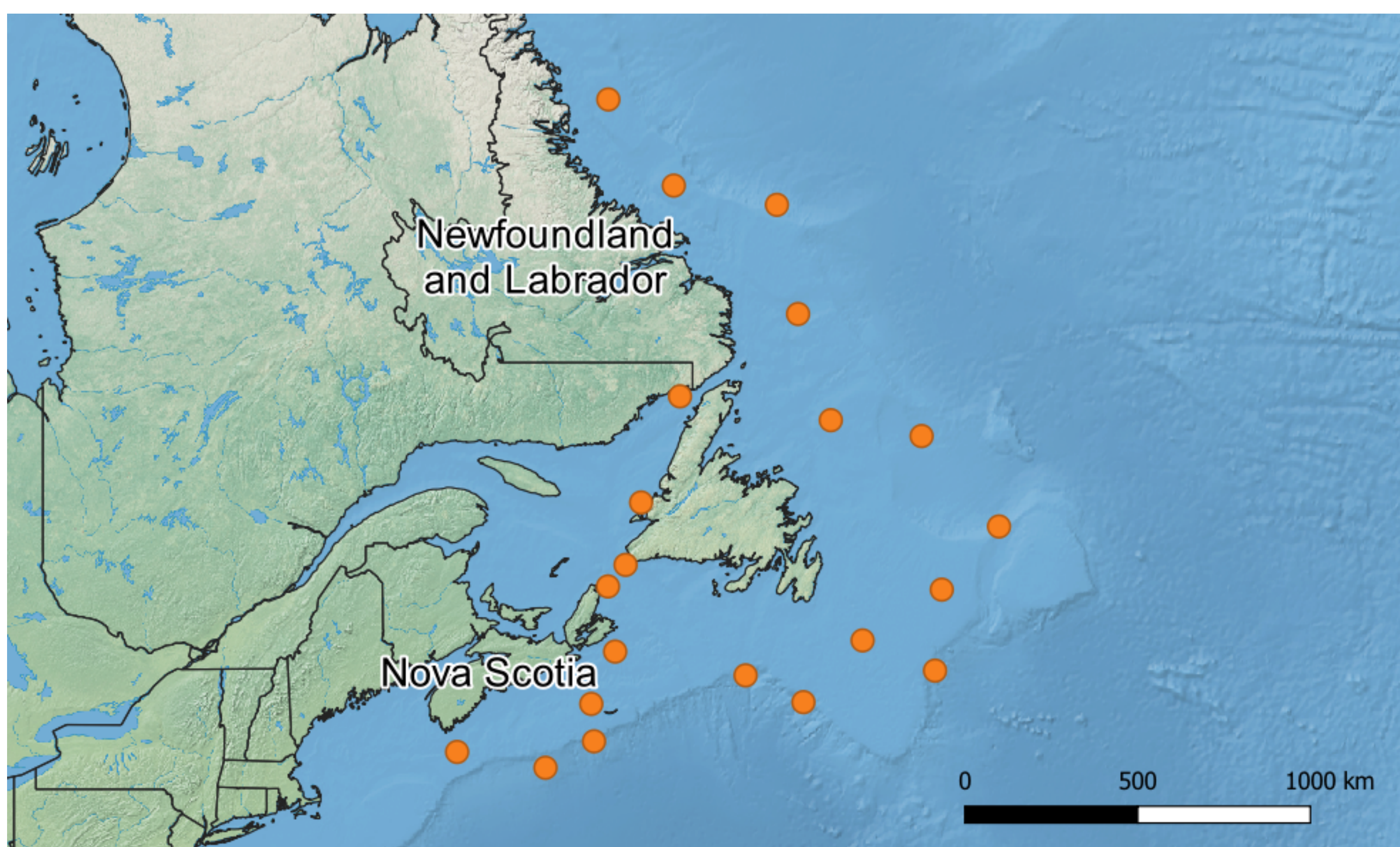
*mark.thomas@dal.ca

Introduction

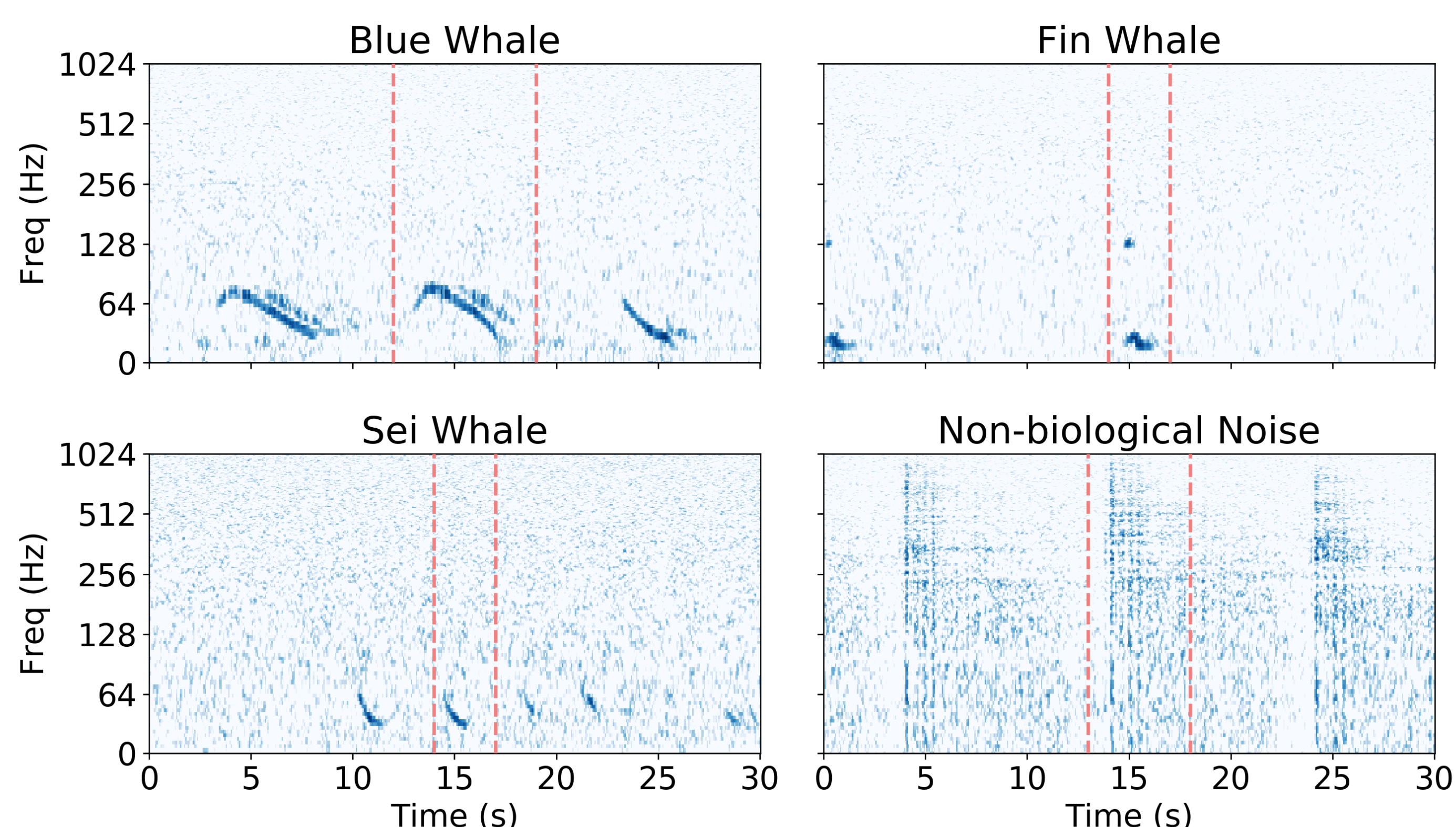
- Marine biologists use acoustic data collected through Passive Acoustic Monitoring (PAM) to determine presence, abundance, behaviour and migratory patterns of marine life, especially marine mammals
- Collections of acoustic recordings obtained through PAM are very large, making complete human analysis infeasible
- Can we use deep learning to detect and classify marine mammal vocalizations in acoustic recordings?

Acoustic Recordings and Training Data

- The acoustic recordings were collected by JASCO Applied Sciences during the summer and fall months of 2015 and 2016 in the areas surrounding the Scotian Shelf



- The recordings were analyzed by marine biologists producing annotations pertaining to marine mammal vocalizations and other acoustic sources labelled as "non-biological"
- We focus on identifying three species of baleen whales with similar call types (blue, fin, and sei whales) against non-biological and ambient sources
- We use spectrograms of the acoustic recordings containing each annotation and treat this problem as an image-classification task



Source	Training	Validation	Test
Blue Whale	2692 (6.23%)	601 (6.49%)	574 (6.20%)
Fin Whale	15118 (35.01%)	3244 (35.06%)	3272 (35.36%)
Sei Whale	1701 (3.94%)	332 (3.59%)	383 (4.14%)
Non-biological	2078 (4.81%)	449 (4.85%)	398 (4.30%)
Ambient	21589 (50.00%)	4626 (50.00%)	4627 (50.00%)

Stacked and Interpolated Spectrograms

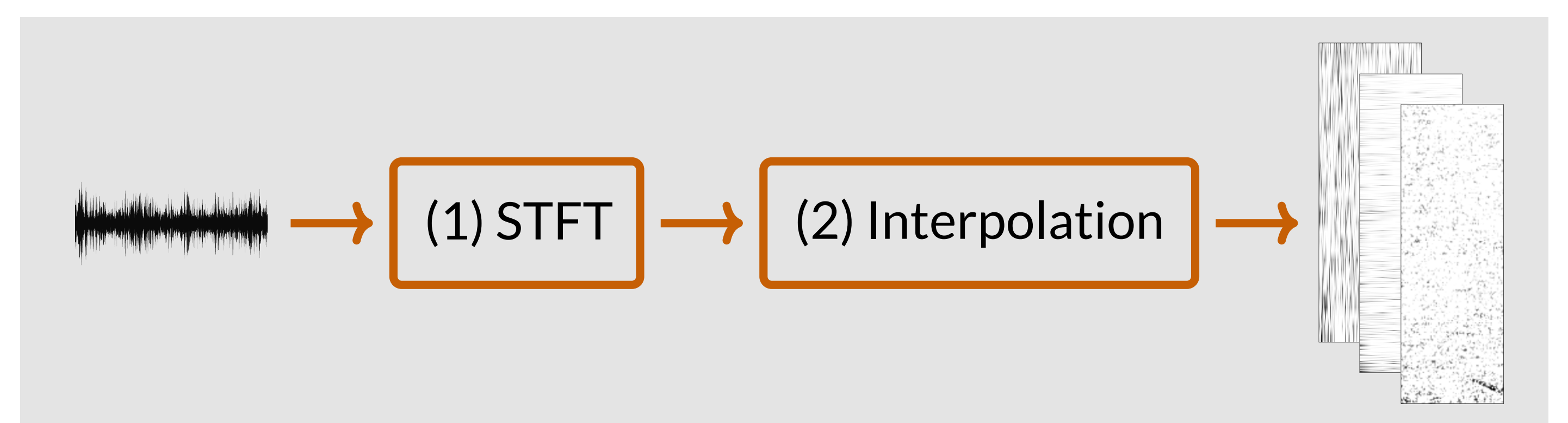
- Experts in marine biology use multiple spectrograms with different resolutions when analyzing acoustic recordings
- How can we exploit the strategy used by marine biologists without simply training multiple classifiers?
 - Generate k spectrograms using multiple sets of parameters to the Short-time Fourier Transform

$$X(n, \omega) = \sum_{m=-\infty}^{\infty} x[m]w[m-n]e^{-j\omega m} \quad (1)$$

- Interpolate the original spectrograms over a pre-defined resolution

$$\omega = \omega_i + \frac{\omega_{i+1} - \omega_i}{n_{i+1} - n_i}(n - n_i) \quad (2)$$

- Stack the interpolated spectrograms to form a k -channel tensor



Neural Network Architecture and Training Details

- We train a commonly used deep Convolutional Neural Network (CNN) known as ResNet-50 [1]
- A cross-entropy loss function was optimized using Stochastic Gradient Descent (SGD) with momentum
- Other training parameters: batch size=128, learning rate=0.001 with exponential decay ($\lambda = 0.01$) every 30 epochs

Experimental Results

	1-channel Standard Spectrogram			3-channel Novel Representation
	NFFT=256	NFFT=2048	NFFT=16384	
Accuracy	0.88512	0.94326	0.94196	0.95331
Precision	0.71979	0.86621	0.85686	0.89265
Recall	0.64634	0.83627	0.83814	0.88409
F-1 Score	0.67394	0.85003	0.84697	0.88735

True label	1-channel Standard Spectrogram					3-channel Novel Representation				
	NFFT=256	NFFT=2048	NFFT=16384			NFFT=256	NFFT=2048	NFFT=16384		
BW	0.27 0.05 0.60 0.06 0.03	0.70 0.04 0.19 0.06 0.01	0.73 0.03 0.16 0.06 0.01	0.75	0.02 0.17 0.06 0.00	0.75	0.02 0.17 0.06 0.00	0.75	0.02 0.17 0.06 0.00	0.75
SW	0.07 0.47 0.42 0.04 0.01	0.03 0.76 0.18 0.03 0.00	0.04 0.75 0.16 0.04 0.01	0.02	0.88 0.09 0.01 0.00	0.02	0.88 0.09 0.01 0.00	0.02	0.88 0.09 0.01 0.00	0.02
FW	0.04 0.01 0.92 0.02 0.01	0.02 0.01 0.95 0.01 0.01	0.02 0.01 0.94 0.01 0.01	0.01	0.96 0.02 0.00	0.01	0.96 0.02 0.00	0.01	0.96 0.02 0.00	0.01
NN	0.10 0.03 0.23 0.58 0.07	0.03 0.02 0.12 0.77 0.06	0.03 0.03 0.13 0.76 0.05	0.04	0.02 0.06 0.85 0.04	0.04	0.02 0.06 0.85 0.04	0.04	0.02 0.06 0.85 0.04	0.04
AB	0.00 0.00 0.00 0.00 1.00	0.00 0.00 0.00 0.00 1.00	0.00 0.00 0.00 0.00 1.00	0.00	0.01 0.00 0.00 0.99	0.00	0.01 0.00 0.00 0.99	0.00	0.01 0.00 0.00 0.99	0.00
	BW SW FW NN AB	BW SW FW NN AB	BW SW FW NN AB		BW SW FW NN AB	BW SW FW NN AB	BW SW FW NN AB		BW SW FW NN AB	
	Predicted label									

References and Acknowledgements

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770--778, 2016.

Collaboration between researchers at JASCO Applied Sciences and Dalhousie University was made possible through an NSERC Engage Grant. The acoustic recordings were collected by JASCO Applied Sciences as part of the Environmental Studies Research Fund (ESRF) program.