

# Segmentation, Classification, and Visualization of Orca Calls using Deep Learning

Hendrik Schröter, Elmar Nöth, Andreas Maier, Rachael Cheng, Volker Barth, **Christian Bergler**  
Pattern Recognition Lab, Friedrich-Alexander University Erlangen-Nürnberg  
IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)  
May 12<sup>th</sup> – 17<sup>th</sup> 2019, Brighton, United Kingdom (UK)



## Motivation: Killer whale research



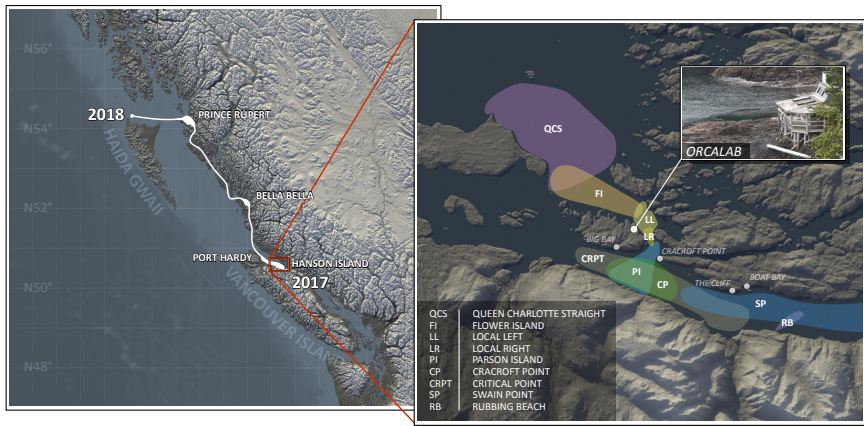
The Killer Whale (*Orcinus orca*) [1]

## Motivation: Killer whale research



OrcaLab [2]

## Motivation: Killer whale research



Covered recording area by the DeepAL [1] expedition and the fixed installed OrcaLab [2] hydrophones



## Motivation: Killer whale research

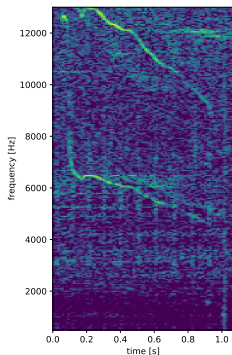
### The Orchive [3]

- collected by the OrcaLab [2] and Steven Ness [3]
- 20,000 hours of underwater recordings by using 6 stationary hydrophones (1985–2010)
- 23,511 digitized audio tapes each  $\sim 45$  min.
- Orchive Annotation Catalog (OAC) [3] comprises 15,480 orca/noise labels

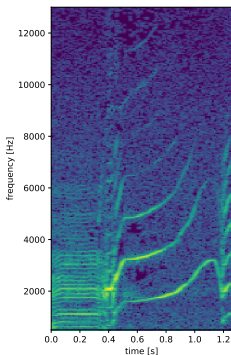
### DeepAL Fieldwork Data (DLFD) 2017/2018 [1]

- collected via a 15-meter research trimaran
- 1,007 hours of multi-channel underwater recordings
- 89 hours video footage about behavioral data

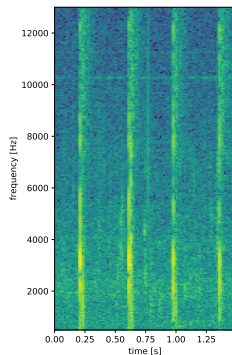
## Example killer whale vocalizations



Whistle



Pulsed Call



Echolocation Click

Spectrograms from three characteristic killer whale sounds.

# Outline

## Data Corpora and Preprocessing

## Segmentation – Network Architecture, Training, and Results

## Call Type Classification – Network Architecture, Training, and Results

## Visualization – Call Type Features

## Conclusion

# Data Corpora and Preprocessing



# Data Corpora – Orca/Noise Segmentation

## Corpora

dataset \ split		training		validation		test	
		samples	% orca	samples	% orca	samples	% orca
OAC <sup>1</sup>	<b>11,504</b>	8,042	<b>84.9</b>	1,711	83.3	1,751	82.4
AEOTD <sup>2</sup>	<b>17,995</b>	14,424	8.9	1,787	15.4	1,784	5.7
DLFD <sup>3</sup>	<b>31,928</b>	23,891	14.2	4,125	30.1	3,912	28.3
SUM	<b>61,427</b>	46,357	<b>24.8</b>	7,623	38.6	7,447	35.6

<sup>1</sup> Orchive Annotation Catalog (OAC) [2]

<sup>2</sup> Automatic Extracted Orchive tape data (AEOTD) [3]

<sup>3</sup> DeepAL Fieldwork Data (DLFD) [1]

# Data Corpora – Call Type Classification

## Corpora

dataset \ split		training		validation		test	
		samples	%	samples	%	samples	%
CCS <sup>1</sup>	<b>138</b>	102	73.9	19	13.8	17	12.3
CCN <sup>2</sup>	<b>286</b>	198	69.2	41	14.4	47	16.4
EXT <sup>3</sup>	<b>90</b>	63	70.0	12	13.3	15	16.7
<b>SUM</b>	<b>514</b>	<b>363</b>	<b>70.6</b>	<b>72</b>	<b>14.0</b>	<b>79</b>	<b>15.4</b>

<sup>1</sup> Call Catalog Symonds (CCS) [2]

<sup>2</sup> Call Catalog Ness (CCS) [3]

<sup>3</sup> Orchive Extension Catalog (EXT)

# Data Preprocessing

## Preprocessing and Augmentation

- Power-Spectrogram
- Augmentation
  - Amplitude scaling
  - Frequency shift
  - Time stretch
  - Addition of noise spectrograms
  - Trimming / Padding to fixed length
- dB-Normalization

# Segmentation – Network Architecture, Training, and Results

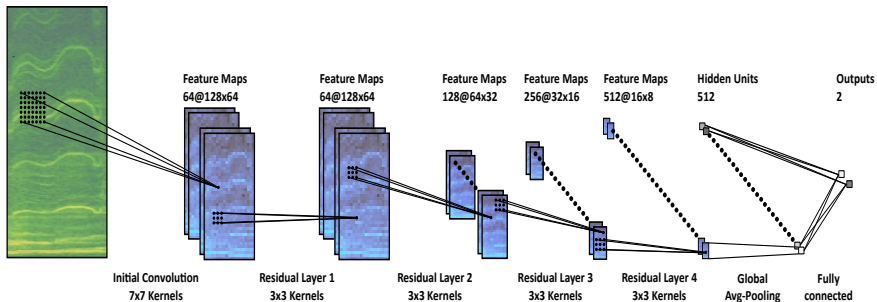




# Network Architecture and Training

## Architecture

Inputs  
1@256x128

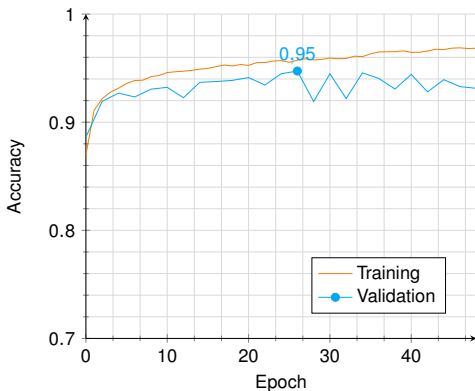


ResNet18-based Convolutional Neural Network (CNN) without max-pooling in the first residual layer for a binary classification problem

## Network Results

### Results

- **Test accuracy of 95.0 %** (TPR = 93.8 %, FPR = 4.3 %)



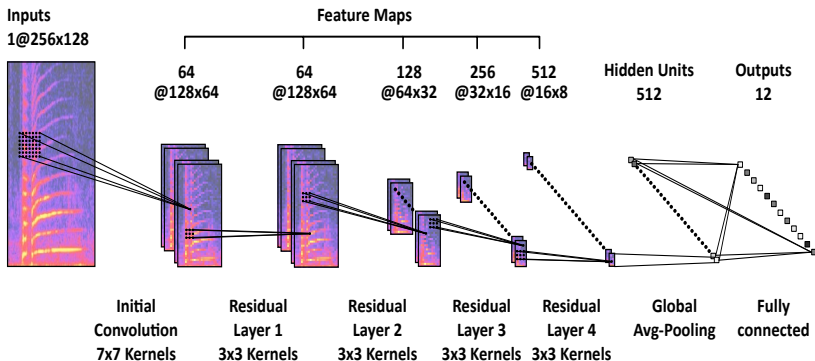
Training and validation accuracy of the segmentation model.

# Call Type Classification – Network Architecture, Training, and Results



# Network Architecture and Training

## Architecture

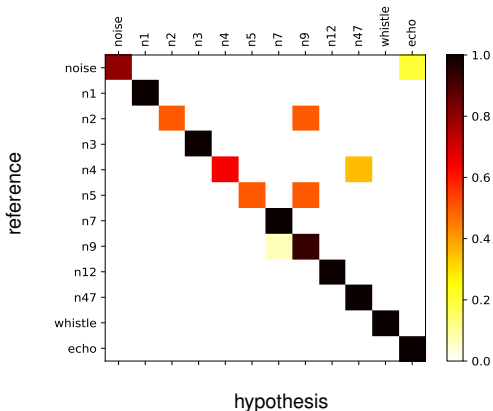


ResNet18-based Convolutional Neural Network (CNN) without max-pooling in the first residual layer for a 12-class problem

## Network Results

### Results

- Mean test accuracy of 87.0%

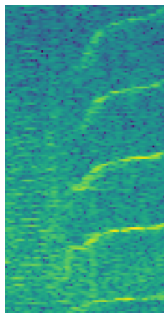


Confusion matrix from the call type classifier.

## Network Results

### Misclassifications

Reference

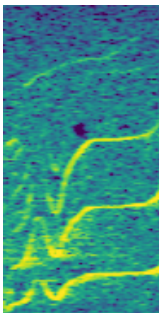


N9

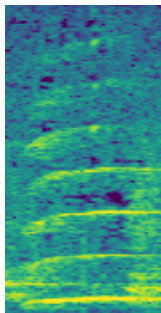
—

Wrong predictions

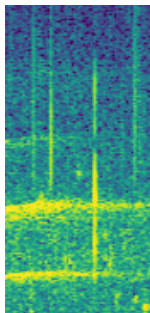
—



N2 as N9



N5 as N9

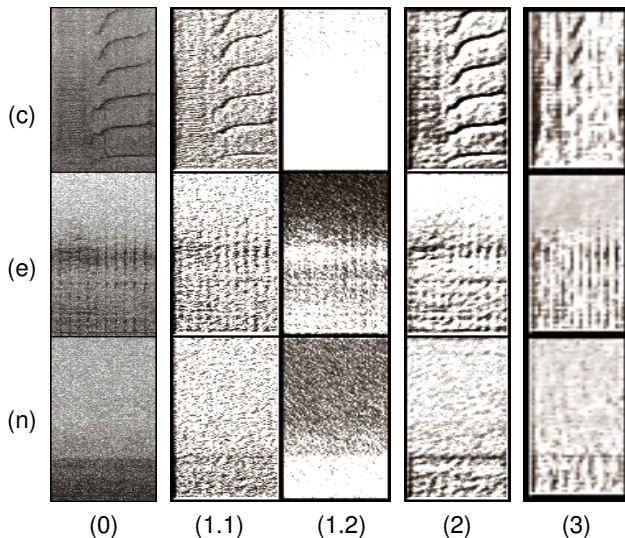


N9 as N7

# Visualization – Call Type Features



## Call Type Feature Visualization





# Conclusion

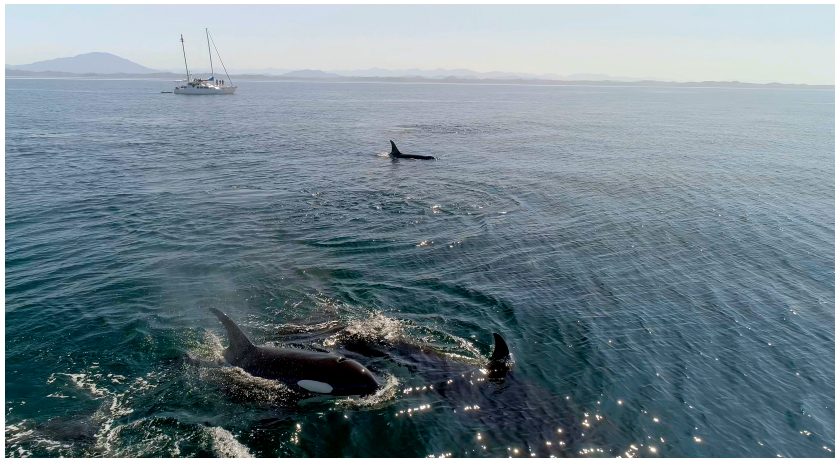


## Conclusion

- Two-stage approach for robust segmentation and classification
- Applicable on any semi-labeled database
- Real-time factor of 1/25 (NVIDIA GTX 1050) enables on-the-fly detection in the field
- Automatically segment large data corpora followed by a subsequent call type classification
- Direct comparison to other work is difficult (different data corpora and/or approaches) (Steven Ness [3])
- Training call type classifier with only few call type labels
- Increase training data to be more robust against signal variety of real-world data

**Thank you for your attention.**

**Questions?**



## References I

- <sup>1</sup>C. Bergler, *Deepal fieldwork data 2017/2018 (dlfd)*, <https://www5.cs.fau.de/research/data/> (April 2019).
- <sup>2</sup>ORCALAB, *Orcalab - a whale research station on hanson island*, <http://orcalab.org> (September 2018).
- <sup>3</sup>S. Ness, “The orchive : a system for semi-automatic annotation and analysis of a large collection of bioacoustic recordings”, PhD thesis (Department of Computer Science, University of Victoria, 3800 Finnerty Road, Victoria, British Columbia, Canada, V8P 5C2, 2013), p. 228.
- <sup>4</sup>A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in pytorch”, in Nips 2017 workshop (2017).

## Data Distribution

### Call Type Label Distribution

Orca Call Type/ Corpus	N01	N02	N03	N04	N05	N07	N09	N12	N47	echo	whistles	noise	SUM
<b>CCS [2]</b>	33	10	—	21	14	18	26	16	—	—	—	—	<b>138</b>
<b>CCN [3]</b>	36	—	56	60	—	31	70	—	33	—	—	—	<b>286</b>
<b>EXT</b>	—	—	—	—	—	—	—	—	—	30	30	30	<b>90</b>
<b>SUM</b>	<b>69</b>	<b>10</b>	<b>56</b>	<b>81</b>	<b>14</b>	<b>49</b>	<b>96</b>	<b>16</b>	<b>33</b>	<b>30</b>	<b>30</b>	<b>30</b>	<b>514</b>

Orca call type, echolocation, whistle, and noise label distribution of the CCS, CCN, and EXT data corpus

# Data Preprocessing

## Preprocessing and Augmentation

**Data:** Training Input Audio  $\mathcal{A}_{inp}$

**Result:** Trainable Spectrogram  $\mathcal{S}_{train}$

- 1  $\mathcal{S}_{inp} \leftarrow 10 \cdot \log_{10}(|\mathcal{FFT}(\text{resamp}(\text{mono}(\mathcal{A}_{inp}), 44.1 \text{ kHz}), \text{ffts} = 4096, \text{hop} = 441)|^2)$
- 2  $\mathcal{S}_{train} \leftarrow \text{scaleAmplitude}(\mathcal{S}_{inp}, \alpha_{dB} = \text{sample}([-6 \text{ dB}, 3 \text{ dB}]))$
- 3  $\mathcal{S}_{train} \leftarrow \text{shiftPitch}(\mathcal{S}_{train}, \alpha = \text{sample}([0.5, 1.5]))$
- 4  $\mathcal{S}_{train} \leftarrow \text{stretchTime}(\mathcal{S}_{train}, \alpha = \text{sample}([0.5, 2]))$
- 5  $\mathcal{S}_{train} \leftarrow \text{compressFrequencies}(\mathcal{S}_{train}, f_{min} = 500\text{Hz}, f_{max} = 10\,000 \text{ Hz}, \text{bins} = 256)$
- 6  $\mathcal{S}_{train} \leftarrow \text{addNoise}(\mathcal{S}_{train}, \text{sample}(\mathcal{S}_{noise}), \text{SNR} = \text{sample}([12 \text{ dB}, -3 \text{ dB}]))$
- 7  $\mathcal{S}_{train} \leftarrow \text{normalize}(\mathcal{S}_{train}, \text{dB}_{min} = -100 \text{ dB}, \text{dB}_{ref} = 20 \text{ dB})$
- 8  $\mathcal{S}_{train} \leftarrow \text{trimPad}(\mathcal{S}_{train}, \text{length} = \text{sample}(128))$
- 9 **return**  $\mathcal{S}_{train}$

## Segmentation Model – Network Training

### Training

- implemented and trained using PyTorch [4]
- Adam optimizer ( $lr_{init} = 10^{-5}$ ,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.999$ )
- learning rate decayed by a factor of 0.5 if there was no improvement on the validation accuracy for 4 epochs
- training stopped if there was no improvement on the validation accuracy for 10 epochs
- batch size = 32

# Classification Model – Network Training

## Training

- implemented and trained using PyTorch [4]
- Adam optimizer ( $lr_{init} = 10^{-5}$ ,  $\beta_1 = 0.5$ ,  $\beta_2 = 0.999$ )
- learning rate decayed by a factor of 0.5 if there was no improvement on the validation accuracy for 4 epochs
- training stopped if there was no improvement on the validation accuracy for 10 epochs
- batch size = 4



## Comparison with previous work: Segmentation

Name	Segment. type	Dataset size	Accuracy	AUC
Ness [3]	Orca	11 041	92.12 %	–
<b>Ours</b>	Orca	61 427	94.97 %	98.17%

## Comparison with previous work: Classification

### Ness [3]

- Classification of 12 pulsed calls
- Mean accuracy of 76 %
- Per class accuracies between 60 % to 92 %

### Ours

- Classification of 9 pulsed calls, whistle, echolocation and noise
- Mean test accuracy of 87 %
- Per class accuracy between 50 % to 100 %