

# ORCA-PARTY: AN AUTOMATIC KILLER WHALE SOUND TYPE SEPARATION TOOLKIT USING DEEP LEARNING

Christian Bergler<sup>1,†</sup>, Manuel Schmitt<sup>1</sup>, Andreas Maier<sup>1</sup>, Rachael Xi Cheng<sup>2</sup>, Volker Barth<sup>3</sup>, Elmar Nöth<sup>1</sup>

<sup>1</sup>Friedrich-Alexander-University Erlangen-Nuremberg, Pattern Recognition Lab, Erlangen, Germany

<sup>2</sup>Leibniz Institute for Zoo and Wildlife Research, Berlin, Germany, <sup>3</sup>Anthro-Media, Berlin, Germany

<sup>†</sup>{christian.bergler}@fau.de

## ABSTRACT

Data-driven and machine-based analysis of massive bioacoustic data collections, in particular acoustic regions containing a substantial number of vocalizations events, is essential and extremely valuable to identify recurring vocal paradigms. However, these acoustic sections are usually characterized by a strong incidence of overlapping vocalization events, a major problem severely affecting subsequent human-/machine-based analysis and interpretation. Robust machine-driven signal separation of species-specific call types is extremely challenging due to missing ground truth data, speaker/source-relevant information, limited knowledge about inter- and intra-call type variations, next to diverse recording conditions. The current study is the first introducing a fully-automated deep signal separation approach for overlapping orca vocalizations, addressing all of the previously mentioned challenges, together with one of the largest bioacoustic data archives recorded on killer whales (*Orcinus Orca*). Incorporating ORCA-PARTY as additional data enhancement step for downstream call type classification demonstrated to be extremely valuable. Besides the proof of cross-domain applicability and consistently promising results on non-overlapping signals, significant improvements were achieved when processing acoustic orca segments comprising a multitude of vocal activities. Apart from auspicious visual inspections, a final numerical evaluation on an unseen dataset proved that about 30 % more known sound patterns could be identified.

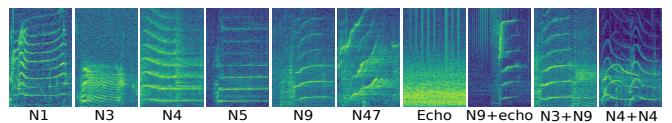
**Index Terms**— Killer Whale, Deep Learning, Call Type

## 1. INTRODUCTION

In the coastal regions of the northeastern Pacific Ocean, more than 40 years have been dedicated to study communication and behavior of the largest member of the dolphin family [1, 2, 3, 4] – the killer whale (*Orcinus Orca*) – resulting in one of the largest animal-specific bioacoustic data repositories. The *Archive* [1, 5, 6, 7] contains 25 years (1985–2010) and approximately 20,000 h of acquired killer whale

Thanks to the German Research Council (DFG) for funding and “The Paul G. Allen Frontier’s Group” for their initial grant. Thanks also to Helena Symonds & Paul Spong (OrcaLab), Steven Ness (formerly UVIC), for giving us permission to use the raw orca-specific data archive, and Simeon Smeele (Max-Planck-Institute), for giving us permission to the Monk parakeet data.

underwater recordings in northern British Columbia while using 6 stationary hydrophones. The vocal repertoire of killer whales consists of three different sound types [3, 8, 7]: (1) *Echolocation Clicks* – short pulses used for navigation and object localization, (2) *Whistles* – narrow-band signals primarily used within close-range interactions, and (3) *Pulsed Calls* – the most common type of vocalizations, subdivided into discrete, variable, and aberrant calls, showing distinct tonal properties. Discrete pulsed calls, also referred to as *Call Types*, are stereotyped and repetitive vocal activities, indicating a wide diversity of distinctive categories with significant inter- and intra-class spectral variations [2, 3, 8, 7]. In order to gain much deeper insights into killer whale communication, it is essential to perform large-scale killer whale call type identification on existing massive bioacoustic data volumes [1, 5, 6] using machine-based, data-driven, and state-of-the-art machine (deep) learning technologies [7]. However, current research on machine-based orca call type recognition [9, 10, 7] is substantially affected by overlapping call type structures (“Cocktail Party Problem” in animal linguistics [11]) in various combinations, such as (1) known vs. known call type, (2) known vs. unknown call pattern, and (3) known call type vs. noise artifacts. Figure 1 visualizes noisy spectral representations of various killer whale vocalizations, including both, non-overlapping but also overlapping samples.



**Fig. 1:** Noisy spectral examples of either different isolated killer whale sound types, or overlapping vocal activities [2]

Large inter- but especially intra-call type variety lead to huge combinatorial variations of possible overlapping paradigms, strongly influencing the high-dimensional feature space. Furthermore, supervised classification also proves to be extremely difficult in case of overlapping structures, since the algorithms are trained exclusively on class-dependent isolated call types [10, 9]. In particular, longer segments with many short successive sequences of vocal interactions, i.e. regions where a lot of vocal activities take place, prove to be extremely valuable for communicative analysis. However, these are exactly the acoustic sections which have

a significantly higher probability of potential overlapping events. In [7], we machine-detected  $\approx 2.5$  million killer whale segments ( $\approx 2,992.02$  h), whereas  $\approx 113,000$  segments were  $> 10$  s ( $\approx 551$  h), which is equal to about 20 % and strongly emphasizing the importance of this work even further. Consequently, a robust, large-scale, data-, and machine-driven method for automatic separation of animal-specific acoustic patterns is of huge importance, but presents numerous challenges: (1) robust machine learning techniques to process massive and noise-heavy data repositories [7], (2) very limited knowledge about the entire inter- and intra call type diversity, (3) huge call type-specific datasets are required to cover as much spectral variation as possible, (4) no ground truth data specifying overlapping calls and the associated individual components, (5) single-channel acoustic events with no information about number of speakers, sound source location, speaker-specific data material, and various recording environments and setups. The current study presents ORCA-PARTY, a fully-automated deep learning-based orca sound type separation framework, independent of speaker-, sound source location-, and recording condition-specific knowledge, not requiring human-annotated overlapping ground truth data.

## 2. RELATED WORK

Sound source separation for human speech is a well-studied field of research which has greatly benefited from the use of various deep learning methods in recent years [12, 13, 14]. In bioacoustics, however, it is comparatively less studied, especially together with modern machine (deep) learning approaches. A transfer of existing human speech-specific sound source separation methods to the field of bioacoustics is non-trivial, due to [15]: (1) wide range of unique animal-specific vocal behaviors (e.g. type of call patterns, call durations, frequency ranges), (2) strong acoustic interference through various environmental conditions, and (3) subjective (human perception) and objective (Scale-Invariant Signal-to-Distortion Ratio (SI-SDR)) qualitative assessment of bioacoustic data compared to human speech. Deng et al. [16] introduced an efficient blind source separation technique, capable of separating signature whistles of dolphins. Hassan et al. [17] applied traditional machine learning techniques (FastICA, PCA, NMF) to separate signal mixtures of overlapping frog sounds. Izadi et al. [18] presents a deep learning-based source separation algorithm, which applied a two-step approach consisting of detection and segmentation to extract bat calls. Zhang et al. [19] developed a bi-directional long short-term memory (BLSTM) network for separation of overlapping bat calls from six different species. Bermant et al. [15] proposed a deep learning-based bioacoustic source separation approach to distinguish between individual-specific vocalizations of 8 rhesus macaque individuals, 8 bottlenose dolphins, and 15 Egyptian fruit bats. To the best of the authors' knowledge, this is the first study addressing deep learning-based orca sound type separation, without requiring any human-annotated overlapping ground truth data, regardless of speaker, sound source location, and recording conditions, to separate overlapping orca signals and enhance downstream call type recognition.

## 3. DATA MATERIAL

### 3.1. Killer Whale Sound Type Archive (KWSTA)

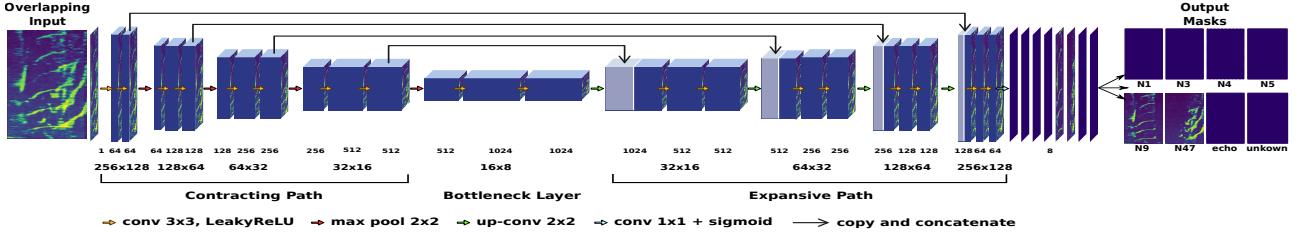
ORCA-SLANG [7] was used for large-scale data-driven identification of orca sounds, processing the entire Orchieve (20,000 h) [1, 6, 5]. The *Killer Whale Sound Type Archive* (KWSTA) includes a total number of 246,852 ( $\approx 398.1$  h) unique orca sound type excerpts (mono, 44.1 kHz sampling rate), consisting of three sub-datasets: (1) *ORCA-SLANG Call Type Data Corpus*, including 235,369 ( $\approx 393.2$  h) audio samples distributed across 6 orca call type categories (N1: 19,280, N3: 8,484, N4: 145,760, N5: 8,861, N9: 41,990, N47: 10,994) with an average duration of  $\approx 6.0$  s, (2) *Echolocation Repository*, containing 9,382 ( $\approx 2.4$  h) echolocation clicks, machine-identified by processing various Orchieve recordings via ORCA-TYPE [9], our ResNet18-based multi-class Convolutional Neural Network (CNN), designed for supervised orca sound type classification. The dataset has an average signal duration of  $\approx 0.9$  s per data sample, and (3) the *ORCA-SLANG Unknown Signal Repository*, comprises 2,101 ( $\approx 2.5$  h) data samples of either so far unseen orca sounds or background noise, with an average signal duration of  $\approx 4.2$  s. Therefor, ORCA-TYPE was applied to classify all remaining ORCA-SLANG [7] data, which could not be assigned to any call type category represented in ORCA-TYPE. Only samples achieving a classification probability of  $< 0.3$  (in case of echos  $< 0.2$ ) were selected. Thus, the probability of observing an unseen killer whale paradigm or noise is very likely.

### 3.2. Call Type Data Corpus (CTDC)

As mentioned in [9, 7] the *Call Type Data Corpus* (CTDC) consists of three individual human-annotated sub-datasets (mono, 44.1 kHz sampling rate): (1) *Orcalab Call Type Catalog* (CCS), (2) *Ness Call Type Catalog* (CCN), and (3) *Extension Catalog* (EXT). Whereas the CCS dataset includes 138 audio files, distributed across 7 different killer whale call types (N1: 33, N2: 10, N4: 21, N5: 14, N7: 18, N9: 26, and N12: 16), the CCN archive comprises 286 call type samples spread across 6 orca call type categories (N1: 36, N3: 56, N4: 60, N7: 31, N9: 70, and N47: 33) [9, 20]. The EXT dataset contains 90 additional audio samples differentiated into echolocation clicks, whistles, and noise, each containing 30 data files. Consequently, the overall CTDC dataset adds up to 514 audio files, categorized into 12 distinct classes [9, 20].

### 3.3. DeepAL Fieldwork Data 2017/2018/2019 (DLFD)

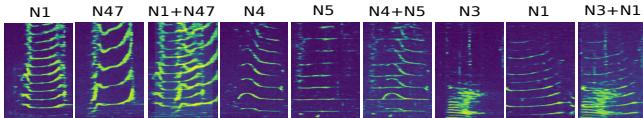
Additional acoustic orca data were collected during our fieldwork expeditions along the coastal waters of northern British Columbia in 2017, 2018, and 2019 [21]. A 15-meter research trimaran was used together with different constellations of a custom-made high sensitivity and low noise towed-array, while following the requirements of the Canadian Department of Fisheries and Oceans [21]. Underwater signals were recorded, digitized, and stored as multi-channel wavfiles at a sampling rate of 96 kHz [21]. The entire *DeepAL fieldwork data 2017/2018/2019* (DLFD) contains  $\approx 177.3$  h (single channel) of additional orca underwater recordings.



**Fig. 3:** ORCA-PARTY – Deep separation network architecture (based on [22]) using a  $256 \times 128$  overlapping input spectrogram (two signals mixed) and  $8 \times 256 \times 128$  network output, modeling all 8 category-specific output masks (Figure modified after [20])

#### 4. METHODOLOGY

To create the signal spectrograms, a multi-stage data pre-processing procedure is applied similar to the proposed workflow in [21, 9, 20]: (1) converting audio signal of various duration to a mono file, (2) re-sampling to 44.1 kHz, (3) STFT (fft-size = 4,096 ( $\approx 100$  ms), hop = 441 ( $\approx 10$  ms)) to create  $F$  (= frequency bins)  $\times T$  (= time frames) decibel-converted power-spectrograms, (4) orca detection algorithm [20], to return a fixed temporal context of 1.28 s, which results in a  $2,049 \times 128$  spectral representation, (5) linear frequency compression (nearest neighbor, fmin = 500 Hz, fmax = 10 kHz) to 256 frequency bins, followed by (6) 0/1-dB-normalization (min =  $-100$  dB, ref = +20 dB) providing the final output of our data preprocessing, a noisy  $256 \times 128$  0/1-dB-normalized spectral representation. In a next step, ORCA-CLEAN was applied to denoise a pair of two  $256 \times 128$  spectrograms, before they were randomly overlapped, using a duration interval  $\delta \in [0.64 \text{ s}, 1.28 \text{ s}]$  leading to an overlap between 50 % up to 100 %. According to the degree of overlap, a temporal context of 1.28 s was randomly sub-sampled and 0/1-min/max-normalized, representing the final  $256 \times 128$ -large input for ORCA-PARTY. To create a combinatorial diverse set of overlapping spectrograms, 5,000 samples for each orca call type category, as well as echolocation clicks, were randomly selected from the KWSTA. In addition, all 2,101 unknown samples of the KWSTA material were selected, leading to a total of 37,101 samples distributed across 8 categories. Based on this data pool, 2,000 overlapping events were randomly generated for each combination (only spectral pairs) resulting in a total of 42 combinations. Thus, the final *ORCA-PARTY Overlapping Dataset* (OPOD) comprises 84,000  $256 \times 128$ -large, 0/1-min/max-normalized, overlapping orca spectrograms, split into: training (58,800 samples, 70.0 %), validation (12,600 samples, 15.0 %), and test (12,600 samples, 15.0 %). Figure 2 visualizes three examples of generated overlapping orca events, consisting of the two  $256 \times 128$  input spectrograms and the overlapping  $256 \times 128$  result. A more detailed description, combined with a guideline for network training and evaluation, is documented here [23].



**Fig. 2:** Spectral input pairs and respective overlapped results

#### 4.1. ORCA-PARTY

ORCA-PARTY is based on the U-Net architecture [22, 20] (see Figure 3). The network is implemented in PyTorch and trained on the OPOD archive. ORCA-PARTY receives as input  $256 \times 128$ -large, 0/1-min/max-normalized, overlapping spectrograms and returns 8 category-specific segmentation masks. An Adam optimizer with an initial learning rate of  $10^{-4}$ ,  $\beta_1 = 0.5$ , and  $\beta_2 = 0.999$  was used, next to the Mean Squared Error (MSE) loss function and a batch size of 32. The best validation loss was considered as evaluation target metric. A learning rate decay of 0.5 after 4 epochs, next to an early stopping criterion after 10 epochs, was used. Max-pooling was used within the contraction path, whereas transposed convolutions were applied for upsampling within the expansive path. All other convolution layers are followed by batch-normalization and LeakyReLU ( $\alpha = 0.1$ ), except the last plain convolutional layer, which is followed by a sigmoid.

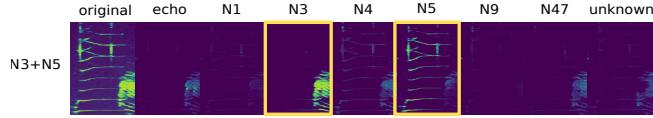
#### 5. EXPERIMENTS

In total three experiments were conducted. First, ORCA-PARTY was trained and verified via a two-stage evaluation process on the unseen OPOD test set. ORCA-PARTY was evaluated by random visual inspection and classification of the final output masks, ignoring overlapping signals involving the “unknown” class, resulting in 8,400 out of 12,600 OPOD test samples. In a second experiment, ORCA-PARTY was deployed as a data enhancement step for supervised call type classification. ORCA-TYPE [9] was trained on the denoised (via ORCA-CLEAN [20]) CTDC material, but only on the 8 categories modelled in ORCA-PARTY, leading to 409 excerpts and a data distribution of: training (289 samples, 70.0 %), validation (58 samples, 15.0 %), and test (62 samples, 15.0 %). The CTDC data does not contain any overlapping signal patterns. A comparison between the supervised trained classifiers – baseline system without (O-BL) and with ORCA-PARTY (O-WP) – was performed as follows: (1) classification performance on human-annotated non-overlapping orca signals, using the CTDC data split, (2) classification performance on machine-detected denoised orca segments (via ORCA-SPOT [21] and ORCA-CLEAN [20]) using our unlabeled/unseen fieldwork recordings (DLFD), while applying a duration interval  $\delta \in [10.0, 30.0]$  seconds. It is assumed that those regions consist of multiple vocal activities and provide a large number of overlapping events (average duration per identified segment in 2017 = 14.8 s, 2018 = 14.7 s, and

2019 = 14.3 s). Both classifiers used a sliding window approach (window = 1.28 s, step-size = 0.64 s) to iterate frame-by-frame over the entire excerpt. O-BL classified each denoised original frame, counting it only if the maximum probability ( $p_{\max}$ ) belonged to one of the 7 sound categories and  $p_{\max} > 0.50$ . O-WP classified the frame-wise content in all 7 masks, excluding the “unknown” category, counting it only if the hypothesis ( $p_{\max} > 0.50$ ) matched the respective type of mask (e.g. N1 classified with  $p_{\max} > 0.50$  in the N1 mask). In case there exist no overlapping segments the number of classifications should be almost identical, while many overlapping structures should result in significantly more classifications of orca sound types, as multiple masks per frame are activated. In a final experiment ORCA-PARTY was transferred to a bird species, named Monk parakeets (*Myiopsitta monachus*).

## 6. RESULTS AND DISCUSSION

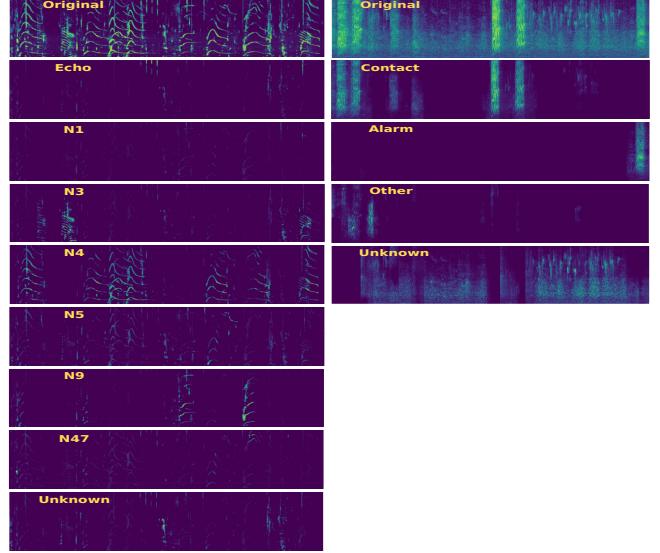
Figure 4 visualizes an example of a randomly-chosen, denoised, overlapping,  $256 \times 128$ -large input spectrogram from the unseen OPOD test set, together with the final category-specific separation outputs, in this case, a multiplication of the original overlapping input spectrogram and the  $8 \times 256 \times 128$  network output. The color-coded categories shown in Figure 4 demonstrate the two strongest class-related activations across all 8 classes, matching the artificially machine-generated overlapping ground truth data. Besides all visual inspections, an overall classification accuracy of  $\approx 86.0\%$  was achieved, while applying O-WP to 8,400 (16,800 output masks/classification hypotheses) unseen OPOD test samples, indicating that in  $\approx 86.0\%$  the best classification hypothesis ( $p_{\max}$ ) matched the category-specific ORCA-PARTY output mask and was part of the ground truth vocalization pair.



**Fig. 4:** Network input and category-based separation output

O-BL, achieved an average accuracy (10 runs, 8 classes) of  $\approx 96.0\%$  (validation) and  $\approx 94.5\%$  (test) with respect to non-overlapping events only, while O-WP obtained an average accuracy of  $\approx 94.5\%$  (validation) and  $\approx 93.0\%$  (test). Due to the small differences in validation and test (both  $\approx 1.5\%$ ), O-WP almost reaches its upper classification boundary, defined by the average classification accuracy of O-BL. In conclusion, non-overlapping orca sound events do not present any problem for ORCA-PARTY. However, the analysis of longer acoustic segments, usually containing substantial amounts of orca vocalizations, significant performance differences between both classifiers appear. O-BL obtained 9,664 vocalization types in 2017, 18,665 in 2018, and 11,240 in 2019 by classifying the entire DLFD archive, summing up to 39,569 orca sound events distributed across 7 classes. O-WP, however, achieved: 12,028 vocalization types in 2017, 24,576 in 2018, and 15,080 in 2019, resulting in 51,684 orca sound types spread across 7 categories. The significant increase of 12,115 ( $\approx 30\%$ ) additional orca sound type hypotheses is

mostly due to the classifier’s inability to interpret overlapping structures often being categorized as noise. Figure 5 (left) visualizes a 17.0 s long denoised orca excerpt of the unseen 2017 DLFD dataset, including plenty of various overlapping orca vocalizations, all of them being correctly identified by ORCA-PARTY and ORCA-TYPE. Besides the original spectral trace, representing the frame-by-frame ORCA-PARTY input, all mask activations are visualized accordingly, following the same rules as in Figure 4.



**Fig. 5:** (left) denoised orca traces (17.0 s), (right) noisy monk parakeet traces (10.0 s), and category-based separation masks

To demonstrate and prove model transferability, we trained ORCA-PARTY on 3,251 human-annotated monk parakeet sound events, divided into 4 classes: 689 contact, 798 alarm, and 764 other calls, next to 1,000 noise files (other songbirds and/or environmental noise). ORCA-PARTY was trained on 3,000 overlapping pair-wise spectral combinations. Figure 5 (right) visualizes a 10 s-long bird excerpt. ORCA-PARTY performs well even in noisy conditions and is able to distinguish across various call categories. It is also noticeable that vocalizations of other songbirds were correctly assigned to the “unknown” class.

## 7. CONCLUSION AND FUTURE WORK

In this study we present ORCA-PARTY, an automatic deep learning-based approach for orca sound type separation. ORCA-PARTY does not require any human-labeled overlapping ground truth data and is independent of speaker-/source information and various recording conditions. As an additional data enhancement step, similar classification results were obtained for non-overlapping events, however, significant improvements were observed during the analysis of acoustic regions with high vocalization volumes, leading to  $\approx 30\%$  more call identifications. ORCA-PARTY has also shown promising initial results on various noisy bird calls. Future studies will evaluate additional animal-related bioacoustic datasets. The source code of ORCA-PARTY will be publicly available under [23], next to audiovisual excerpts.

## 8. REFERENCES

- [1] Steven R. Ness, *The Archive : A system for semi-automatic annotation and analysis of a large collection of bioacoustic recordings*, Ph.D. thesis, Department of Computer Science, University of Victoria, British Columbia, Canada, 2013.
- [2] John K. B. Ford, “A catalogue of underwater calls produced by killer whales (*Orcinus orca*) in British Columbia,” *Canadian Data Report of Fisheries and Aquatic Science*, , no. 633, pp. 165, 1987.
- [3] John K. B. Ford, “Vocal traditions among resident killer whales (*Orcinus orca*) in coastal waters of British Columbia,” *Canadian Journal of Zoology*, vol. 69, pp. 1454–1483, June 1991.
- [4] J.K.B. Ford, G.M. Ellis, and K.C. Balcomb, *Killer whales: The natural history and genealogy of Orcinus orca in British Columbia and Washington*, UBC Press, 2000.
- [5] Steven R. Ness, “Orchive,” <http://orchive.cs.uvic.ca/> (October 2021).
- [6] ORCALAB, “Orcalab - A whale research station on Hanson Island,” <http://orcalab.org> (October 2021).
- [7] Christian Bergler, Manuel Schmitt, Andreas Maier, Helena Symonds, Paul Spong, Steven R. Ness, George Tzanetakis, and Elmar Nöth, “ORCA-SLANG: An Automatic Multi-Stage Semi-Supervised Deep Learning Framework for Large-Scale Killer Whale Call Type Identification,” in *Proc. Interspeech*, 2021.
- [8] John K. B. Ford, “Acoustic behaviour of resident killer whales (*Orcinus orca*) off Vancouver Island, British Columbia,” *Canadian Journal of Zoology*, vol. 67, pp. 727–745, January 1989.
- [9] Christian Bergler, Manuel Schmitt, Rachael Xi Cheng, Hendrik Schröter, Andreas Maier, Volker Barth, Michael Weber, and Elmar Nöth, “Deep Representation Learning for Orca Call Type Classification,” in *Proc. Text, Speech, and Dialogue 2019*. 2019, vol. 11697 LNAI, pp. 274–286, Springer.
- [10] Christian Bergler, Manuel Schmitt, Rachael Xi Cheng, Andreas Maier, Volker Barth, and Elmar Nöth, “Deep Learning for Orca Call Type Identification – A Fully Unsupervised Approach,” in *Proc. Interspeech*, 2019.
- [11] Mark Bee and Christophe Micheyl, “The cocktail party problem: What is it? how can it be solved? and why should animal behaviorists study it?”, *Journal of comparative psychology (Washington, D.C. : 1983)*, vol. 122, pp. 235–51, 09 2008.
- [12] DeLiang Wang and Jitong Chen, “Supervised speech separation based on deep learning: An overview,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. PP, 08 2017.
- [13] Yi Luo and Nima Mesgarani, “Conv-TasNet: Surpassing ideal time-frequency magnitude masking for speech separation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, pp. 1256–1266, 2019.
- [14] Neil Zeghidour and David Grangier, “Wavesplit: End-to-end speech separation by speaker clustering,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2840–2849, 2021.
- [15] Peter C. Bermant, “Biocppnet: Automatic bioacoustic source separation with deep neural networks,” *bioRxiv*, 2021.
- [16] Xiaohong Deng, Yi Tao, Xingbin Tu, and Xiaomei Xu, “The separation of overlapped dolphin signature whistle based on blind source separation,” in *2017 IEEE International Conference on Signal Processing, Communications and Computing (ICSPCC)*, 2017, pp. 1–5.
- [17] Norsalina Hassan and Dzati Ramli, “A comparative study of blind source separation for bioacoustics sounds based on fastica, pca and nmf,” *Procedia Computer Science*, vol. 126, pp. 363–372, 01 2018.
- [18] Mohammad Rasool Izadi, Robert Stevenson, and Laura N. Kloepfer, “Separation of overlapping sources in bioacoustic mixtures,” *The Journal of the Acoustical Society of America*, vol. 147, no. 3, pp. 1688–1696, 2020.
- [19] Kangkang Zhang, Liu Tong, Shengjing Song, Zhao Xin, Shijun Sun, Walter Metzner, Jiang Feng, and Ying Liu, “Separating overlapping bat calls with a bi-directional long short-term memory network,” *Integrative zoology*, 04 2021.
- [20] Christian Bergler, Manuel Schmitt, Andreas Maier, Simeon Smeele, Volker Barth, and Elmar Nöth, “ORCA-CLEAN: A Deep Denoising Toolkit for Killer Whale Communication,” in *Proc. Interspeech*, 2020.
- [21] Christian Bergler, Hendrik Schröter, Rachael Xi Cheng, Volker Barth, Michael Weber, Elmar Nöth, Heribert Hofer, and Andreas Maier, “Orca-spot: An automatic killer whale sound detection toolkit using deep learning,” *Scientific Reports*, vol. 9, 12 2019.
- [22] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. 2015, vol. 9351 of *LNCS*, pp. 234–241, Springer.
- [23] Christian Bergler, “Open Source Github-Repository,” <https://github.com/ChristianBergler>.