

# Advancing Marine Bioacoustics with Deep Generative Models: A Hybrid Augmentation Strategy for Southern Resident Killer Whale Detection

Bruno Padovese<sup>a,\*</sup>, Fabio Frazao<sup>a,b</sup>, Michael Dowd<sup>c</sup>, Ruth Joy<sup>a</sup>

<sup>a</sup>*School of Environmental Science, Simon Fraser University, Burnaby, V5A 1S6, BC, Canada*

<sup>b</sup>*Faculty of Computer Science, Dalhousie University, Halifax, B3H 1W5, NS, Canada*

<sup>c</sup>*Department of Mathematics and Statistics, Dalhousie University, Halifax, B3H 1W5, NS, Canada*

## Abstract

Automated detection and classification of marine mammals vocalizations is critical for conservation and management efforts but is hindered by limited annotated datasets and the acoustic complexity of real-world marine environments. Data augmentation has proven to be an effective strategy to address this limitation by increasing dataset diversity and improving model generalization without requiring additional field data. However, most augmentation techniques used to date rely on effective but relatively simple transformations, leaving open the question of whether deep generative models can provide additional benefits. In this study, we evaluate the potential of deep generative for data augmentation in marine mammal call detection including: Variational Autoencoders, Generative Adversarial Networks, and Denoising Diffusion Probabilistic Models. Using Southern Resident Killer Whale (*Orcinus orca*) vocalizations from two long-term hydrophone deployments in the Salish Sea, we compare these approaches against traditional augmentation methods such as time-shifting and vocalization masking. While all generative approaches improved classification performance relative to the baseline, diffusion-based augmentation yielded the highest recall (0.87) and overall F1-score (0.75). A hybrid strategy combining generative-based synthesis with traditional methods achieved the best overall performance with an F1-score of 0.81. We hope this study encourages further exploration of deep generative models as complementary augmentation strategies to advance acoustic monitoring of threatened marine mammal populations.

**Keywords:** Deep Learning, Underwater Bioacoustics, Southern Resident Killer Whales

## 1. Introduction

Killer Whales (*Orcinus orca*) are a widespread species inhabiting diverse marine environments across all oceans. The species comprises several distinct ecotypes, some of which can be further divided into populations with each adapted to specific ecological niches and characterized by unique behaviors and intricate social structures. Their extensive vocal repertoire spans frequencies from 0.5 to 100 kHz and includes three primary call types: echolocation clicks, single-toned whistles, and discrete pulsed calls (Ford et al., 1987; Ford and Fisher, 1978), each serving distinct behavioral functions. Echolocation clicks, concentrated between 20 and 100 kHz, are primarily used for navigation and prey detection (Barrett-Lennard et al., 1996; Au et al., 2004). Whistles (0.5–25 kHz) are primarily used for social interactions within pods (Miller, 2006; Thomsen et al., 2001; Riesch et al., 2008). Pulsed calls, which occur in the same frequency range, are the most common vocalization type, supporting group cohesion, individual identification, and more complex social communication (Ford et al., 1987; Ford and Fisher, 1978; Filatova et al., 2009). Furthermore, population-specific vocal dialects add a layer of complexity to their acoustic repertoire, reflecting the social and cultural divergence that distinguishes populations within and across ecotypes.

In the Northeast Pacific, Killer Whales have evolved into three genetically and culturally distinct ecotypes (Barrett-Lennard et al., 1996; Morin et al., 2024), which frequently share overlapping habitats: Resident, Transient, and Offshore Killer Whales (Ford et al., 1998; Riesch et al., 2012; Baird and Stacey, 1988). The Resident (fish-eating) ecotype contains two populations, the Northern Resident population that is made up of three linguistic clans, and the Southern Resident population which consists of one linguistic clan. The populations remain reproductively isolated (Morin et al., 2024; Riesch et al., 2012), and there are strong cultural and linguistic barriers between clans. The Southern Resident Killer Whale population (SRKW) also known as J-clan, is made up of three stable, matrilineal family groups, identified as J, K and L pod. The SRKW range stretches from California to southeast Alaska, and with heavy use in the cross-boundary waters of the Salish Sea. The population is listed as endangered in Canada and is protected by the Species at Risk Act (SARA) <sup>1</sup> and by the Endangered Species Act <sup>2</sup> in the US. By the end of 2024, The Center for Whale Research, the organization responsible for semiannual SRKW census numbers, estimated there were only 75 individuals remaining <sup>3</sup>.

\*Corresponding author

Email address: bruno\_padovese@sfu.ca (Bruno Padovese)

<sup>1</sup><https://species-registry.canada.ca/index-en.html#/species/699-5>

<sup>2</sup><https://www.fisheries.noaa.gov/topic/laws-policies/endangered-species-act>

<sup>3</sup>The Center for Whale Research Report

In light of this critical population status, there is significant interest from researchers, citizen scientists, government agencies, and First Nation groups in the protection and conservation of SRKW and their habitat. In many cases, visual sightings remain the primary method of monitoring these individuals (Olson et al., 2018; Sato et al., 2021), but this approach is constrained by weather conditions and the need for active, daytime surface observations, which are predominantly conducted by citizen scientists. More recently, the availability of affordable acoustic recording devices and expanded data storage capacity has enabled large-scale passive bioacoustic monitoring, leading to the creation of extensive audio datasets (Roch et al., 2017; Dede et al., 2014; Webster and Budney, 2017). However, a lack of highly specialized person-hours for bioacoustic data analysis creates a significant bottleneck in the manual detection, annotation, and validation of whale calls across large acoustic datasets. This is a time-consuming process that can span weeks or months for a single hydrophone deployment. Consequently, there is a need for efficient processing methods that can automate key steps in the workflow, such as vocalization classification and detection, within a reasonable timeframe (Stowell, 2022) of a few hours or days instead of weeks. This demand for automated solutions has driven the development of SRKW detection and classification algorithms based on signal processing techniques (Gillespie et al., 2013; Shapiro and Wang, 2009) and machine learning (ML) (Sharpe et al., 2019; Brown et al., 2010). Within ML, deep learning (DL) approaches, originally developed for applications in image, speech, and music processing (Abeßer, 2020; Manilow et al., 2020) have gained traction due to their success in complex pattern recognition and feature extraction (LeCun et al., 2015; Goodfellow et al., 2016). Deep Neural Networks (DNN) have been shown to outperform traditional ML techniques (Stowell, 2022; Morfi et al., 2021), including in the context of Killer Whale acoustic detection and classification (Bergler et al., 2019; Hauer et al., 2023; Bergler et al., 2021).

Two-dimensional spectrograms, typically segmented into fixed-length audio clips (e.g., 1-second or 10-second duration), are commonly used as inputs to DNNs for cetacean classification and detection tasks (Kirsebom et al., 2020; Bergler et al., 2019; Li et al., 2020). This practice is also commonplace across other fields such as sound event classification (Ozer et al., 2018), bird song recognition (Kahl et al., 2021), as well as speech (Wang and Chen, 2018) and music classification (Elbir and Aydin, 2020). In marine bioacoustics, spectrograms are also one of the primary tools for visualizing and analyzing acoustic data, allowing researchers to quickly identify patterns that may be missed through manual listening, especially for sounds outside the human audible range. Furthermore, many acoustic signals, such as Killer Whale vocalizations, contain frequency-modulated (FM) components that are discernible in spectrograms (Rabiner and Juang, 1993), making them suitable for use in automated classification models.

However, while spectrogram-based DNNs have shown promise (Kirsebom et al., 2020; Bergler et al., 2019; Shiu et al., 2020), their effectiveness is constrained by the availability of sufficient high-quality annotated data (Priestley et al., 2023;

Gudivada et al., 2017), a common issue in marine bioacoustics. As mentioned above, obtaining high-quality annotations is challenging and expensive, particularly for marine environments where target vocalizations can be sparsely distributed against a backdrop of overwhelming underwater environmental noise (Bergler et al., 2019; Stowell, 2022; Padovese et al., 2023). This scarcity of annotated data represents a critical bottleneck in the development of effective DL models for marine mammal classification. Therefore, before investing in costly and time-consuming efforts to manually annotate and assign labels to acoustic data, an effective strategy to mitigate this limitation is to use data augmentation to artificially enhance the diversity of small bioacoustic training dataset (Stowell, 2022; Li et al., 2021).

Data augmentation has long been used to address challenges associated with small or unbalanced datasets by creating modified versions of existing data (Shorten and Khoshgof-taar, 2019). These augmented samples are variations of the original recordings that were not present in the training set but are theoretically possible within the same context. Generally, augmentation methods can be divided into two categories: “naïve” augmentations and data-based augmentations. Naïve augmentations are characterized by their simplicity and do not take the specific problem being addressed into consideration. In bioacoustics, commonly used naïve methods include time-shifting (Shiu et al., 2020), time and frequency masking (Park et al., 2019), and noise addition (Mishachandar and Vairamuthu, 2021). These methods are popular due to their ease of implementation, effectiveness in improving model performance, and “safety”, meaning they are unlikely to compromise the meaning of the audio, ensuring the augmented signals remain consistent with the original class (Stowell, 2022).

In contrast, data-based augmentations leverage domain knowledge to transform samples based on the characteristics of the environment and target species. These can include simple operations like sound mixing (Padovese et al., 2021) or more complex techniques like sound propagation modeling, which simulates how sound is distorted as it propagates through water (Binder, 2018). Some studies have experimented with more risky transformations, including warping (Park et al., 2019), pitch shifting (Li et al., 2021), and time stretching (Li et al., 2021). While these methods can potentially increase the diversity of the training data, they also risk distorting subtle acoustic features that are unique to specific species or call types (Stowell, 2022). As a result, the appropriate choice of augmentation techniques is highly context-dependent and should be tailored to the characteristics of each dataset and species.

While traditional data augmentation helps address the limitations of small or unbalanced datasets, it is inherently constrained by the nature of the transformations themselves, which can only recombine or distort existing information in limited ways. To further expand the dataset and introduce new variability, data synthesis that relies on algorithms to generate artificial data has emerged as a valuable tool (King et al., 2014; Reichert and Ronacher, 2015). Already widely used in marine mammal communication research (King et al., 2014; Reichert and Ronacher, 2015), synthetic data can also serve as artificial

training examples to supplement or even replace real-world data in machine learning models (Li et al., 2020). Importantly, synthetic data generation and classical data augmentation strategies are not mutually exclusive; combining both can yield superior results in deep learning applications by maximizing the diversity of training data.

Within data synthesis, DL-based generative models learn the distribution of the dataset’s feature space to create new, unseen samples drawn from this distribution, introducing entirely new patterns not found in the original dataset. Generative models have long been used for tasks such as speech and music generation (Shorten and Khoshgoftaar, 2019). Some studies have begun applying generative models in the field of bioacoustics; this field is still emerging, with limited research having made its way into the peer-review literature. Most research has focused on the generation aspect, typically using Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) or Variational Autoencoders (VAEs) (Kingma, 2013) to operate on time-frequency spectrograms (Zhang et al., 2022; Nieto-Mora et al., 2024). Although a handful of studies have also applied generative models specifically for data augmentation (Herbst et al., 2024; Li et al., 2023), even fewer have systematically compared these methods to traditional data augmentation approaches.

Recently, denoising diffusion probabilistic models (DDPMs) (Ho et al., 2020) have gained attention for their ability to generate higher-quality samples compared to their GAN or VAE counterparts (Dhariwal and Nichol, 2021a). Unlike these earlier approaches, DDPMs generate high-quality, diverse samples by gradually refining noise into informative data, offering better stability compared to traditional GANs or VAEs. These models have shown significant potential for producing realistic synthetic data and have been studied for improving neural network training in tasks such as image classification (Trabucco et al., 2023). In bioacoustics, diffusion models remain largely unexplored, with, to the best of our knowledge, few studies investigating their potential for this purpose (Herbst et al., 2024).

In this work, an aim is to extend the application of DDPMs for data augmentation tailored to SRKW vocalizations. Our approach assumes that, similar to traditional augmentation, additional latent information can be derived from the original dataset through the use of deep generative models. By combining generative models with simpler augmentation techniques, we aim to improve the overall performance of deep learning classifiers, while overcoming the limitations of each method when used independently, particularly under conditions of scarce annotated data.

Our work makes three key contributions: (i) we introduce the first application of diffusion models for data augmentation of SRKW vocalizations; (ii) we propose a hybrid approach integrating diffusion-based synthetic data generation with traditional augmentation techniques to enhance dataset diversity and model robustness; and (iii) we conduct a comparative evaluation of generative versus traditional data augmentation methods, addressing a gap in the literature by systematically assessing their impact on classification performance.

## 2. Materials and Methods

### 2.1. Datasets

The underwater acoustic data used in this study were collected from two hydrophone deployments in the Salish Sea, off the west coast of North America (Figure 1). We focused on discrete pulsed calls from SRKW pods J, K, and L, with all vocalizations manually verified as true positives. Non-tonal sounds, such as whistles and echolocation clicks, were excluded. The recordings were obtained using different systems and protocols, each with distinct site characteristics, as detailed below.

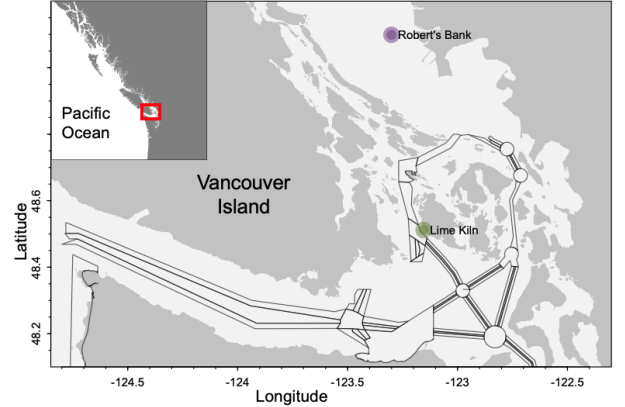


Figure 1: Locations of the two hydrophone deployments in the Salish Sea (Lime Kiln and Roberts Bank). The commercial shipping lanes are shown as black lines.

The first dataset was collected near Lime Kiln State Park, off San Juan Island ( $48^{\circ}30'42''$  N,  $123^{\circ}09'15''$  W), through deployments conducted by SMRU Consulting Ltd. This dataset includes 1,633 audio files, each lasting one minute, recorded between August 29, 2018, and October 16, 2019. The recordings were captured using a hydrophone deployed at a depth of approximately 23 meters and sampled at 250 kHz. To identify Killer Whale vocalizations, recordings were first processed using the PAMGuard whistle and moan detector (Gillespie et al., 2008), which generated initial binary detections indicating the presence or absence of potential biological sounds. All detected files were then manually reviewed in full, and vocalizations from various marine mammal species and ecotypes were annotated. From these, only annotations corresponding to confirmed KW calls were retained for this study, resulting in a total of 1,261 annotations from 200 files. This dataset was used exclusively for model training and hyperparameter tuning.

The second dataset was collected at Robert’s Bank, British Columbia ( $49^{\circ}01'07.8''$  N,  $123^{\circ}11'32.8''$  W), by JASCO Applied Sciences. It comprises 1,562 5 minute audio files, recorded between September 21, 2015, and April 12, 2018. Recordings were made using a hydrophone deployed at a depth of approximately 168 meters and sampled at 64 kHz. Killer Whale encounters were initially identified using a proprietary detection algorithm developed by JASCO Applied Sciences. These encounters were then manually reviewed and annotated by expert analysts for the presence of Killer Whale vocalizations. For this study, we selected only those annotations con-

firmed to originate from KWs, resulting in 1,263 annotations across 22 files. These files and annotations were used exclusively for testing.

The annotations were made publicly available through the HALLO (Humans and Algorithms Listening and Looking for Orcas) project. The full set of annotations can be accessed via the project’s GitHub repository.<sup>4</sup> The underlying audio data were originally released in Palmer et al. (2025).

### 2.1.1. Data Preparation

Killer Whale audio segments were extracted from recordings according to the annotations. Each labeled segment was isolated as a 3-second audio clip, a duration long enough to capture most SRKW calls (Frazao et al., 2025), while short enough to avoid overwhelming the neural network with excessive background noise. For the background class, which contains only (non-KW) background noise, random segments were drawn from the recordings while avoiding overlap with annotated regions. All recordings were downsampled to 24 kHz with anti-aliasing to ensure consistent processing across datasets with original sampling rates ranging from 64-250 kHz. This cutoff frequency was chosen as the resulting Nyquist frequency (12 kHz) fully encompasses the fundamental frequencies and harmonics characteristic of SRKW calls (Ford, 1989) while avoiding unnecessary computational overhead from higher sampling rates. To create a balanced dataset, the number of background segments randomly selected was made equal to the number of SRKW clips, and both were combined to form the complete baseline training dataset.

### 2.1.2. Spectrogram Computation

We computed 128-band Mel spectrograms from the 3-second segments derived from the annotation windows (Section 2.1.1). Spectrograms were generated using a 50 ms Hann window (NFFT of 1200 samples) with a 12.5 ms hop length (300 samples), representing a 75% overlap at the 24 kHz sampling rate. This configuration produced Mel spectrograms with frequency coverage up to 12 kHz, adopting a similar parameterization that has been successfully employed in SRKW classification tasks (Frazao et al., 2025). Amplitude values were then converted to a decibel scale for subsequent neural network processing.

### 2.2. Time-Shifting Augmentation

Time-shifting artificially expands the dataset by temporally displacing audio clips within their original recordings. For a given spectrogram  $x$  of duration  $T$ , we generate  $N$  augmented samples by shifting the annotation window forward and backward in increments of 0.5 seconds. Each shifted window maintains a minimum overlap of 50% with the original annotation, ensuring that the shifted segments remain contextually relevant to the labeled content without introducing mismatched examples. This guarantees that every annotation is represented by at least  $N \geq 1$  instances, though the exact value of  $N$  depends on the original annotation duration and the overlap requirement.

Beyond simply increasing sample size, time-shifting improves model robustness to natural temporal variations in audio events. Unlike methods such as noise addition (White et al., 2022) or random masking (Park et al., 2019), it preserves the original acoustic features while maximizing the value of the limited labeled data.

### 2.3. Vocalization Mask Augmentation

While most augmentation strategies target within-class variability, such as differences in pitch, duration, or timing, they often overlook the variability introduced by natural soundscapes. This contextual information, such as various sources of transient sounds and shifting ambient noise conditions, can influence how a DNN distinguishes vocalizations from background events. Ignoring this context can limit model generalization, particularly in real-world deployments where acoustic conditions vary widely. To address this and create realistic, context-aware synthetic vocalizations, we developed a vocalization mask augmentation strategy based on high-quality SRKW call examples.

To create the masks, we used a curated catalogue of high-quality SRKW vocalizations as the source for clean signal references. The samples in this collection were compiled over several decades of research by Dr. John Ford and made publicly available<sup>5</sup>. We first computed spectrograms from the recordings following the procedure described in Section 2.1.1. We then projected the spectrograms into a lower-dimensional space using PCA. In this representation, the first principal component tends to capture the broad scale pattern that is consistent across time and frequency, and explain the most variance in the dataset. These typically correspond to persistent background noise and non-whale acoustic features present in the environment. By subtracting this component from each original spectrogram, we can obtain a mostly denoised representation that emphasized the vocalization while suppressing background content.

Next, to further refine these masks, we applied a thresholding step in which all pixel values below the  $i$ -th percentile of the spectrogram’s dynamic range were set to zero. This process suppressed residual background noise resulting in sparse, high-contrast masks that captured most vocal features. The full process of mask construction and refinement is illustrated in Figure 2.

The resulting masks were then linearly combined with randomly sampled background spectrograms from the Robert’s Bank dataset. Specifically, we selected segments from recordings that were outside the 22 files reserved for the test set and manually verified to contain only background noise. By combining the masks with actual environmental noise, we created new, high-fidelity vocalizations that reflect the acoustic complexity encountered in the field. This augmentation approach exposes the model to more realistic combinations of signal and background, ultimately improving its robustness and generalization to unseen acoustic environments.

<sup>4</sup><https://github.com/coastal-science/hallo-data/tree/main>

<sup>5</sup><https://orca.research.sfu.ca/call-library/home.html?v=20240530-1727>



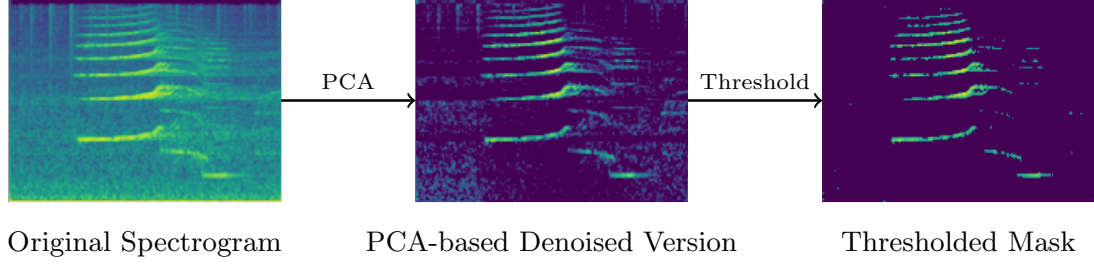


Figure 2: Overview of the vocalization mask construction process. Starting from the original spectrogram (left), we first apply PCA-based background subtraction to emphasize vocalization components (center). A subsequent percentile-based thresholding step produces a sparse, high-contrast mask (right) that preserves the primary vocal features while suppressing residual background noise.

In practice, however, the manual validation of background segments described above may not even be necessary. In real-world applications, randomly sampling unlabeled audio from long-term passive acoustic recordings is likely to overwhelmingly yield only background noise. This makes the approach both scalable and easy to implement.

#### 2.4. VAEs

Variational Autoencoders (VAEs) (Kingma, 2013) are a class of generative models that learn to generate new samples by encoding inputs into compact representations and decoding them back into the original data format. This encoder-decoder structure allows the model to capture the essential features of the input data. The intermediate representation, commonly referred to as the latent space, serves as a compressed version of the data, which VAEs learn to organize in a smooth and continuous manner suitable for generation.

In a standard autoencoder, the encoder learns to compress an input  $x$  into a lower-dimensional representation  $z$  through successive layers, gradually compressing the information into a simplified form that captures its most important features. The decoder does the opposite, attempting to reconstruct the input as  $\hat{x} = p_\theta(x | z)$ , where  $\hat{x}$  is the reconstructed input, and  $p_\theta(x | z)$  represents the decoder network. This network gradually expands the compact representation back into a full-resolution output. Crucially, the decoder can only learn to accurately reconstruct inputs if the encoded representation is informative and semantically meaningful. To achieve this, the model is trained to minimize the difference between the input and its reconstruction, encouraging it to preserve essential structure while discarding noise or redundancy. However, because the encoder only learns to handle inputs it has seen during training, the resulting latent space can be irregular and fragmented, even if it compresses data well. As a result, there’s no guarantee that randomly sampling from this space will produce valid or meaningful outputs, making standard autoencoders poorly suited for generative tasks.

VAEs address this by introducing a probabilistic encoding scheme. Instead of mapping each input to a single fixed point in the latent space, the encoder  $q_\phi(z | x)$  learns to represent it as a Gaussian distribution defined by a mean vector  $\mu$  and a standard deviation vector  $\sigma$ . A sample is drawn from this normal distribution using the so-called reparameterization trick:

$$z = \mu + \sigma \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (1)$$

where  $I$  denotes the identity covariance matrix, following the standard VAE formulation.

The decoder  $p_\theta(x | z)$  then reconstructs the input from this sampled vector. Unlike in standard autoencoders, where  $z$  might be sampled from uninformative parts of the latent space, VAEs are explicitly trained to make their latent space well-organized and continuous. This is done by nudging the learned distributions  $q_\phi(z | x)$  to stay close to a known prior, such as a standard normal distribution  $\mathcal{N}(0, I)$ . As a result, VAEs learn to fill the latent space smoothly, such that small changes in  $z$  correspond to gradual changes in the decoded output. This structured organization makes it possible to sample new, realistic data simply by drawing  $z \sim \mathcal{N}(0, I)$  and passing it through the decoder.

The model is trained to minimize the following loss (Kingma, 2013):

$$\mathcal{L}_{\text{VAE}} = -\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x | z)] + \beta \text{KL}(q_\phi(z | x) || p(z)) \quad (2)$$

where the first term is the reconstruction loss, encouraging the decoder to accurately reconstruct the input  $x$  from the sampled latent representation  $z$ , and the second term is a Kullback-Leibler (KL) divergence that regularizes the encoder’s output distribution to remain close to the standard normal prior  $p(z) = \mathcal{N}(0, I)$ .  $\beta$  is a weight that controls the strength of the regularization; in this work, it was set to 1.

The stability and ability of VAEs to capture a wide range of data distributions make them advantageous for data augmentation tasks (Shorten and Khoshgoftaar, 2019). However, a common drawback of VAEs is that their outputs tend to be less sharp or detailed than the original inputs or compared to those produced by other generative models, such as GANs (Kingma et al., 2019; Wang et al., 2020). This blurriness arises from the probabilistic nature of the decoder and the averaging effect of the reconstruction loss, which can smooth out fine-grained features.

In this work, we employed a standard convolutional VAE architecture in which both the encoder and decoder are composed of stacked convolutional layers with batch normalization and ReLU activations. The model takes as its input 2D spectrograms. The generated spectrograms aim to match the overall

distribution of the Lime Kiln training data. Further implementation details, including architectural specifications and training configuration, are provided in Section 2.9.

## 2.5. GANs

Generative Adversarial Networks (GANs) are another category of generative models that learn to synthesize data by jointly training two competing neural networks: a generator  $G$ , which produces synthetic samples intended to resemble those from the training distribution, and a discriminator  $D$ , which attempts to distinguish between real and generated samples (Goodfellow et al., 2014). The central idea behind GANs is that, as the discriminator improves at identifying fake samples, the generator must produce increasingly realistic outputs in order to fool it. Conversely, as the generator becomes better at producing convincing samples, the discriminator must also improve to maintain its ability to detect fakes. This dynamic coupling forms an adversarial feedback loop in which both networks iteratively enhance their capabilities. During training, the generator and discriminator are optimized simultaneously in a zero-sum game, where each network’s objective directly opposes the other. Once the adversarial training process stabilizes, the generator can be used independently to synthesize new samples, whether images, or, in our case, spectrograms of SRKW vocalizations.

GANs have been widely applied in artificial image generation (Karras et al., 2019; Choi et al., 2018), particularly in domains such as facial image synthesis (Karaoglu et al., 2021), scene reconstruction (Wang et al., 2018), and image-to-image translation (Isola et al., 2017). While their application to time–frequency representations such as spectrograms is less common, some works have demonstrated the potential of GANs in the bioacoustic domain (Li et al., 2023; Bergler et al., 2022; Shim et al., 2021). Unlike natural image synthesis, spectrogram generation of vocalizations involves fine, curvilinear features such as SRKW harmonic ridges, which are more sensitive to any distortion.

Formally, the generator network maps a random noise vector  $z \sim \mathcal{N}(0, I)$  to the data space, producing a synthetic spectrogram  $\hat{x} = G(z)$ . The discriminator  $D(x)$  receives either a real spectrogram  $x \in X_{\text{real}}$ , drawn from the data distribution  $p_{\text{data}}$ , or a synthetic one  $\hat{x}$ , and outputs a scalar representing the probability that the input is real. The networks are trained with opposing objectives defined by the minimax loss function:

$$\min_G \max_D \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(\hat{x}))]. \quad (3)$$

While GANs have demonstrated impressive capabilities in synthesizing visually realistic images, they present several limitations. First, GANs are notoriously difficult to train and highly sensitive to hyperparameter choices, network architecture, and optimization dynamics (Arjovsky and Bottou, 2017). These factors often lead to unstable training, non-convergent behavior, or poor-quality outputs (Agarwal and Farid, 2021). Second, GANs are prone to mode collapse (Che et al., 2016), a failure in which the generator produces a limited set of outputs, such as

repeatedly generating highly similar or identical samples (Srivastava et al., 2017; Mariani et al., 2018; Dhariwal and Nichol, 2021b). This leads to a lower diversity of generated data compared to the real distribution, limiting the model’s ability to represent the full variability of the training data, which is an important drawback when using GANs for data augmentation when training a classifier.

Following Goodfellow et al. (2014), a number of advancements have focused on improving training stability and sample quality. These include architectural refinements such as the Deep Convolutional GAN (DCGAN) (Radford et al., 2015), as well as alternative training objectives like the Wasserstein GAN (WGAN) (Arjovsky et al., 2017) and its gradient-penalized variant, WGAN-GP (Gulrajani et al., 2017). Ultimately, however, despite these improvements, GANs still suffer from fundamental training challenges and instability.

In this work, we adopted the DCGAN architecture, which introduces architectural constraints such as convolutional layers without fully connected components, batch normalization, and ReLU/LeakyReLU activations to improve training stability and sample quality. Our implementation follows the original GAN formulation described above, with both the generator and discriminator operating on 2D spectrogram tensors. The generator produces time–frequency representations that aim to match the distribution found in the Lime Kiln training set. The model was trained using the standard GAN loss. A detailed description of the network architecture and training routine is provided in Section 2.9.

## 2.6. DDPM

Denosing Diffusion Probabilistic Models (DDPMs) are a class of generative models that synthesize data through an iterative denoising process (Ho et al., 2020). DDPMs operate in two phases: a forward process, which incrementally corrupts training samples (e.g., spectrograms) by adding Gaussian noise over discrete timesteps, and a reverse process, which learns to iteratively recover the original data by predicting and removing this noise (Nichol and Dhariwal, 2021). During training, the model learns the underlying structure of the data distribution by estimating the noise at each corruption step (Sohl-Dickstein et al., 2015). During inference, the trained model “hallucinates” novel samples (spectrograms in our case) by progressively denoising pure random noise. This reverse trajectory can generate realistic outputs that closely approximate the statistical distribution of the training dataset, enabling the creation of entirely new, yet plausible, synthetic vocalizations (Kong et al., 2020; Herbst et al., 2024).

The forward process incrementally transforms a clean SRKW spectrogram  $x_0 = x$  into pure noise  $x_T$  over  $T$  timesteps (i.e., discrete diffusion steps) using a predefined noise schedule (Figure 3). The noise schedule dictates how the noise is gradually added at each timestep through the noise schedule coefficients  $\alpha_t$ . These coefficients  $\alpha_t \in (0, 1)$  determine the proportion of the original signal retained at each timestep, with smaller values of  $\alpha_t$  introducing more noise. The cumulative product of these coefficients,  $A_t = \prod_{s=1}^t \alpha_s$ , determines the amount of the

original clean signal  $x_0$  that is preserved at timestep  $t$ . At each timestep  $t \in (1, T)$ , the noisy sample  $x_t$  is computed as:

$$x_t = \sqrt{A_t} x_0 + \sqrt{1 - A_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (4)$$

where  $\epsilon$  is a noise sample drawn from a standard Gaussian distribution (analogous to  $z$  in the VAE formulation above).

The reverse process aims to recover the clean spectrogram  $x_0$  from the noise sample  $x_T$  by progressively denoising over  $T$  timesteps. A neural network, typically a U-Net (Ronneberger et al., 2015), is trained to predict the added noise at each timestep. The U-Net architecture is particularly well-suited for this task due to its ability to capture hierarchical features at different resolutions. During training, the network predicts the noise  $\epsilon$  at each timestep, minimizing the simplified ( $\mathcal{L}_{\text{simple}}$ ) loss function:

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t, x_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2] \quad (5)$$

where  $\mathcal{L}_{\text{simple}}$  represents the mean-squared error (MSE) between the actual noise  $\epsilon$  and the noise predicted by the network  $\epsilon_\theta(x_t, t)$ .

As the model is trained to minimize the MSE loss function across many noisy samples, it gradually becomes better at removing the noise step by step, encouraging the network to become better at removing noise from the spectrogram and reconstructing the clean spectrogram,  $x_0$ . Once trained, the model synthesizes novel spectrograms by iteratively denoising random noise over  $T$  timesteps.

DDPMs overcome many limitations of GANs and VAEs. Rather than relying on adversarial training or latent reconstruction, they minimize a denoising objective at each timestep, leading to more stable training and higher sample quality and diversity (Cao et al., 2024). As a result, DDPMs are capable of generating highly detailed samples that often surpass those produced by both GANs and VAEs (Ho et al., 2020; Dhariwal and Nichol, 2021b).

## 2.7. Filtering Low-Quality Synthetic Spectrograms

Despite the success of DDPMs, and generative models in general, in producing visually appealing and realistic images, several limitations remain, particularly when such models are used for data augmentation in classification tasks. Generative models, including GANs and DDPMs, may synthesize samples containing artifacts or failure regions that can negatively impact downstream classifier performance. This issue can be especially pronounced in the context of marine bioacoustics, where spectrograms often contain low signal-to-noise ratios and complex background interference (e.g., vessel noise or other environmental sounds). In such cases, generative models may inadvertently learn to replicate the noise rather than the vocalization signal of interest. Figure 4 showcases examples of common failures in synthetic spectrograms generated from noisy marine mammal datasets.

To mitigate the risk of low-quality or out-of-distribution synthetic samples negatively affecting the classifier model training, we implemented a simple PCA-based filtering strategy de-

signed to align the statistical distribution of generated spectrograms with that of real SRKW data.

Let  $X_{\text{real}}$  and  $X_{\text{gen}}$  denote the sets of real and generated spectrograms, respectively. To assess the quality of synthetic data, we applied PCA to  $X_{\text{real}}$  to define a low-dimensional projection space, and projected both  $X_{\text{real}}$  and  $X_{\text{gen}}$  into this space. We then used the Mahalanobis distance to quantify how closely the projection of a generated spectrogram  $z_{\text{gen}}$  aligned with the distribution of  $X_{\text{real}}$  in this space.

A generated spectrogram was retained if its Mahalanobis distance  $d_M$  satisfied:

$$d_M(z_{\text{gen}}) \leq \tau \quad (6)$$

where  $\tau$  is a threshold set to the  $j$ -th percentile of Mahalanobis distances computed from  $X_{\text{real}}$ .

We applied this filtering procedure to the output of all generative models, to ensure that only synthetic samples statistically aligned with the distribution of real SRKW vocalizations were used for data augmentation.

## 2.8. Deep Learning Classifier

Our aim is to generate synthetic vocalizations using the augmentation methods described in the previous sections, and to use these data to train a DNN to detect SRKW calls in novel acoustic environments within a prescribed level of accuracy. To this end, we trained a convolutional neural network (CNN) to classify spectrograms as containing SRKW vocalizations, or not.

We employed a ResNet-18 architecture (He et al., 2016), a compact variant of the residual network family. Given the limited size of our training dataset the ResNet-18 strikes a balance between representational capacity and computational efficiency. The network is designed to accept 3-second Mel-spectrogram representations as input, and to output a binary classification probability indicating the presence or absence of SRKW vocalizations.

While the ResNet-18 was the architecture of choice in this study, numerous deep learning classifiers have demonstrated success in detecting and classifying marine mammal sounds. Classical CNN architectures, such as LeNet and VGG-16, as well as GRU RNNs, have shown robust performance in classifying North Atlantic Right Whale (NARW) vocalizations in early deep learning applications in marine bioacoustics (Shiu et al., 2020). More recent architectures, such as DenseNets and Inception models have further advanced performance in related tasks (Tiwari et al., 2023). Indeed, the ResNet architecture itself has been widely adopted in the field (Padovese et al., 2021; Murphy et al., 2022), including for Killer Whale classification off the west coast of Canada (Bergler et al., 2019). While our focus was not on identifying the optimal classifier, it is reasonable to expect that alternative architectures could also provide satisfactory performance based on these prior successes. Ultimately, the central challenge in bioacoustics remains addressing the constraints posed by complex noise environments and small annotated datasets rather than selecting the perfect neural network architecture.

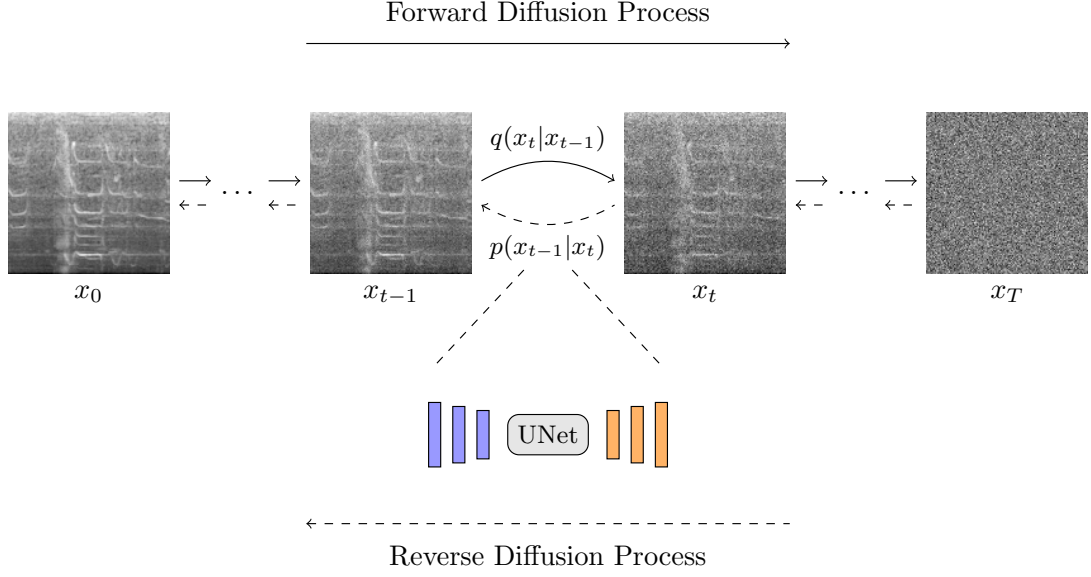


Figure 3: Illustration of a Denoising Diffusion Probabilistic Model (DDPM) for synthesizing spectrograms. The forward diffusion process (solid arrow) incrementally corrupts a clean input spectrogram  $x_0$  into pure noise  $x_T$  over  $T$  timesteps. The reverse diffusion process (dashed arrow) learns to denoise by training a U-Net to predict and remove noise at each step.

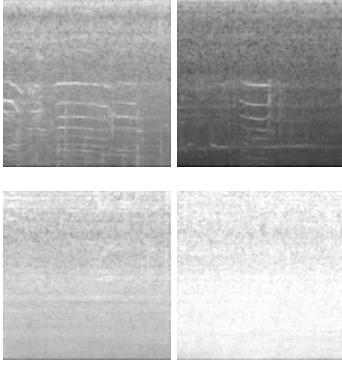


Figure 4: Examples of synthetic spectrograms generated using DDPMs. The top row shows samples that were accepted for training, exhibiting clear SRKW-like vocal structure. The bottom row shows samples that are unsuitable due to the presence of artifacts or poor signal definition.

## 2.9. Experimental Setup

To evaluate the effectiveness of the data augmentation strategies, we conducted a series of experiments comparing classifier performance across seven training regimes: (I) baseline (no augmentation), (II) time-shifting augmentation only, (III) vocalization mask augmentation only, (IV) VAE-generated synthetic data only, (V) GAN-generated synthetic data only, (VI) DDPM-generated synthetic data only, and (VII) hybrid augmentation combining time-shifting, vocalization masks, and the best-performing generative model, selected empirically based on evaluation performance (see Section 3). Models were trained using vocalizations from the Lime Kiln dataset and evaluated on vocalizations from the independent Robert’s Bank dataset, allowing us to assess generalization to a site not used for training.

All experiments were carried out on a dedicated workstation

equipped with an NVIDIA GeForce GTX 1070ti GPU (8 GB memory). In the following, we will describe the model architectures and parameters used for training the different models.

**Classifier.** Each classifier was trained using the same ResNet-18 architecture described in Section 2.8, trained for 20 epochs with a batch size of 128 samples. The model was optimized using an Adam optimizer with a learning rate of 0.001, and a cosine annealing scheduler with linear warmup. As the generative models operated on  $128 \times 128$  spectrograms, all input data were resized to  $128 \times 128$  when necessary, and normalized to  $[0, 1]$ .

**VAE.** The VAE model consisted of a convolutional encoder with three layers (32, 64, 128 filters, kernel size 4, stride 2), each followed by ReLU activation, and two fully connected layers to compute the mean and log-variance of a 32-dimensional latent vector. The decoder mirrored this structure with a fully connected layer followed by three transposed convolutional layers (128, 64, 32 filters), each with ReLU activation, and a final output layer with tanh activation. Models were trained for 150 epochs using a batch size of 64, and a learning rate of 0.001 with an Adam optimizer. The loss combined MSE reconstruction and KL divergence. Input images were resized to  $128 \times 128$  pixels and normalized to  $[-1, 1]$ .

**GANs.** The GAN model followed a DCGAN-style architecture, with a generator composed of six transposed convolutional layers (with 1024, 512, 256, 128, 64, and 1 filters, respectively), each followed by BatchNorm and ReLU activations, except for the final layer which used tanh. The discriminator consisted of five convolutional layers with LeakyReLU activations and PhaseShuffle layers (Donahue et al., 2018) to promote invariance to local time shifts. Batch normalization was applied after



all but the first convolutional layer. Both networks were trained using the Adam optimizer with a learning rate of 0.0001 and  $(\beta_1, \beta_2) = (0.5, 0.999)$ , for 300 epochs with a batch size of 64. Input images were resized to  $128 \times 128$  pixels and normalized to the  $[-1, 1]$  range.

**DDPM.** For diffusion-based synthesis we adopted a standard DDPM built around a U-Net backbone. The U-Net processes  $128 \times 128$  single-channel spectrograms and comprises six resolution levels, with two residual blocks per level. The encoder used six blocks with 128, 128, 256, 256, 512, and 512 filters, respectively. The fifth block, operating at  $8 \times 8$  resolution, incorporated self-attention to better capture long-range time–frequency structure, while the remaining blocks used standard downsampling operations. The decoder mirrored this structure in reverse.

We used the original DDPM formulation with 1,000 diffusion steps and a cosine  $\beta_t$  schedule, training for 150 epochs with a batch size of 32. Optimization employed an AdamW optimizer (learning rate  $1 \times 10^{-4}$ ,  $\beta_1=0.9$ ,  $\beta_2=0.999$ ) and an exponential moving average of the network weights (decay 0.9999) for evaluation. During training the objective was the MSE between predicted and true noise. All inputs were resized to  $128 \times 128$  pixels and normalized to the  $[-1, 1]$  range, matching the other generative models for consistent downstream comparison.

**Experimental Protocol.** The baseline model in experiment I was trained using only real data, that is, 1,261 3-second SRKW vocalization clips extracted from annotated regions and an equal number of background segments sampled from non-annotated intervals (Section 2.1.1). This model establishes a performance reference point against which all augmented variants were compared.

For Experiments II through VII, augmented vocalization samples were generated independently and added to the baseline training set to create expanded datasets. In each case, we augmented the number of vocalizations by two different amounts, +5,000 and +20,000 samples, to investigate how classifier performance is affected by the quantity of synthetic data introduced during training. For Experiment VII, the hybrid strategy combined time-shifting, vocalization masking, and the best-performing generative model. To keep the total number of augmented samples consistent with the other experiments, the same augmentation budget of +5,000 and +20,000 new samples was used. This total was divided equally among the three augmentation methods, resulting in +1,667 new samples per method for the +5,000 setting and 6,667 per method for the +20,000 setting. A summary of the experimental combinations and vocalization sample sizes is presented in Table 1.

To maintain class balance across all training sets, each set of vocalization samples, whether real or augmented, was paired with an equal number of background samples. In the baseline and time-shifting experiments (I and II), background segments were randomly drawn from non-annotated intervals in the Lime Kiln recordings. For the remaining experiments (III–VII), background samples were generated using the same strategy

as their corresponding vocalizations. In the case of vocalization masking (Experiment III), since the masks were overlaid onto background spectrograms from the Robert’s Bank dataset (see Section 2.3), we randomly sampled additional background segments from Robert’s Bank recordings that were outside the recordings reserved for the test set.

For all synthetic augmentation methods (Experiments IV–VII), a separate generative model of the same type as the one used for vocalizations was trained exclusively on background data from Lime Kiln. For example, background samples in Experiment IV were generated using a VAE trained on background spectrograms, while Experiment V used a background GAN, and so forth. This approach ensured that both positive and negative samples in each experiment were drawn from comparable distributions and generated under consistent modeling assumptions. In light of this, throughout the text, when we refer to experiments being augmented by 5,000 or 20,000 vocalization samples, it should be understood that an equal number of background samples was also added, following the corresponding method for that experiment.

Each experiment was repeated 10 times with different random initializations. Performance was quantified using precision, recall, and F1-score, with final metrics averaged across all runs to establish the mean performance. The code used in this study has been made publicly available at <https://github.com/bpadovese/GAugSRKW>, accompanied by comprehensive documentation and a command-line interface to facilitate reuse and adaptation for other bioacoustics applications.

### 3. Results

We illustrate a visual comparison between 12 real SRKW vocalization spectrograms in Figure 5 (a), and the same number of synthetic samples generated by VAE (b), GAN (c), and DDPM (d) respectively. This subset is intended as a representative but small sample illustrating the qualitative differences among models. Visually, we observed that the VAE-generated samples suffered from a characteristic over-smoothing effect. Harmonic ridges appear blurred, and in many cases, signal boundaries are poorly defined or diffused into the background.

The GAN samples, by contrast, displayed sharper structures and better-defined signals than the VAE outputs. However, despite avoiding the most severe cases of mode collapse, the GAN still seemed to produce a limited range of variation. Several spectrograms appear overly similar in composition, with vocalizations often occurring in the same time–frequency locations.

Finally, the DDPM synthetic spectrograms exhibited high visual fidelity, preserving fine harmonic structure and modulation contours that are characteristic of discrete pulsed calls, while avoiding the over-smoothed appearance often observed with VAEs and the low sample variety typical of GANs.

Table 2 summarizes the classification performance in terms of precision, recall, and F1-score at a decision threshold of 0.5, a commonly adopted default in classification tasks used in many studies (Herbst et al., 2024; Duc, 2020; Li et al., 2023). Time-shifting augmentation (II) led to a substantial increase in recall,

Table 1: Summary of experiment combinations and number of vocalization samples. A dash indicates that no samples from that augmentation strategy were included. For Experiment VII, only the best-performing generative model was used. Each set of vocalizations was paired with an equal number of background samples generated using the same method.

Exp.	Real	Time-Shift	Masks	VAE	GAN	DDPM
I ( <i>Baseline</i> )	1,261	–	–	–	–	–
II ( <i>Time-Shift</i> )	1,261	+5,000 +20,000	–	–	–	–
III ( <i>Masks</i> )	1,261	–	+5,000 +20,000	–	–	–
IV ( <i>VAE</i> )	1,261	–	–	+5,000 +20,000	–	–
V ( <i>GANs</i> )	1,261	–	–	–	+5,000 +20,000	–
VI ( <i>DDPM</i> )	1,261	–	–	–	–	+5,000 +20,000
VII ( <i>Hybrid</i> )	1,261	+1,667 +6,667	+1,667 +6,667	–	–	+1,667 +6,667

reaching 0.62 and 0.71 for the +5,000 and +20,000 settings, respectively, compared to just 0.42 for the baseline model. Precision remained comparable to the baseline at +5,000 samples (0.70 vs 0.65), but improved notably to 0.76 when the number of augmented samples increased to 20,000, resulting in F1-scores of 0.66 and 0.73, respectively. The mask-only strategy (III) consistently achieved near-perfect precision (0.98 and 0.99) at both augmentation levels, with overall comparable performance across settings. Recall remained relatively modest, rising slightly from 0.51 at +5,000 to 0.53 at +20,000.

The VAE-based augmentation (IV) did not provide substantial gains at the +5,000 augmentation level, with performance remaining comparable to the baseline (F1-score of 0.51). In contrast, the +20,000 setting resulted in a modest improvement, increasing the F1-score to 0.57. Similarly, GAN-based augmentation (V) produced comparable results, achieving F1-scores of 0.52 and 0.60. In contrast, the DDPM-based approach (VI) demonstrated consistently stronger performance, with F1-scores of 0.71 and 0.75, representing an improvement of approximately 0.15 over the other generative models at both augmentation levels. Notably, it achieved the highest recall overall, reaching 0.87 in the +20,000 setting. Finally, the combined augmentation approach (VII) delivered the best overall performance in the +20,000 setting, combining high precision at 0.99, with a recall of 0.69, resulting in the highest F1-score of 0.81.

Next, we analyzed the classification performance of the proposed methodology across all seven experimental regimes under different decision thresholds. Figure 6 displays the mean precision-recall curves for all seven experiments described in Table 1, showing classification performance with respect to the SRKW class for the +20,000 augmentation setting. Each curve represents the average over ten training runs. The gray line corresponds to the baseline model (I), trained exclusively on real vocalizations. The blue and orange lines represent time-shifting (II) and vocalization mask (III) augmentations, respectively. The green, red, and purple lines correspond to generative augmentation strategies using VAE (IV), GAN (V), and DDPM

(VI). Finally, the brown line (VII) shows the performance of the combined approach using time-shifting, masking, and DDPM-generated samples.

Unlike Table 2, which reports precision, recall, and F1-score at a fixed decision threshold of 0.5, the precision–recall curves in Figure 6 illustrate performance across the entire range of thresholds, and provide a more complete picture of model behavior. The baseline model (I) achieved the lowest performance across the board, with steep precision decay at low recall levels, reflecting the limited generalization when trained only on a small set of real samples from Lime Kiln. As expected, applying any augmentation method improved performance to varying degrees. Traditional augmentation methods maintained precision above 0.9 for recall values below 0.5, with the vocalization mask model (III) outperforming time-shifting (II) across most of the curve. Notably, the mask-only model achieved near-perfect precision at the cost of recall, showing a sharp decline as recall increased above 0.6.

Among the generative approaches, VAE-based (IV) and GAN-based (V) augmentation produced similar precision–recall profiles, both outperforming the baseline but lagging behind traditional augmentation strategies by more than  $\sim 0.1$  in precision across most of the curve. In contrast, the DDPM model (VI) consistently outperformed not only the VAE and GAN models but also the time-shifting approach (II) across the entire recall range. We note, however, that the DDPM model only surpassed the vocalization masking model (III) at the upper end of the recall spectrum (above  $\sim 0.8$ ), where it maintained higher precision.

Finally, the best overall performance was achieved by the combined strategy (VII), which consistently maintained the highest precision for any recall. The improvement was especially pronounced at higher recall values, with near-perfect precision remaining stable up to recall  $\sim 0.7$  and only dropping below 0.9 beyond  $\sim 0.8$ . Even at these higher recall levels, the model preserved a margin of  $\sim 0.2$  in precision compared to all other strategies, highlighting the improved generalization to the

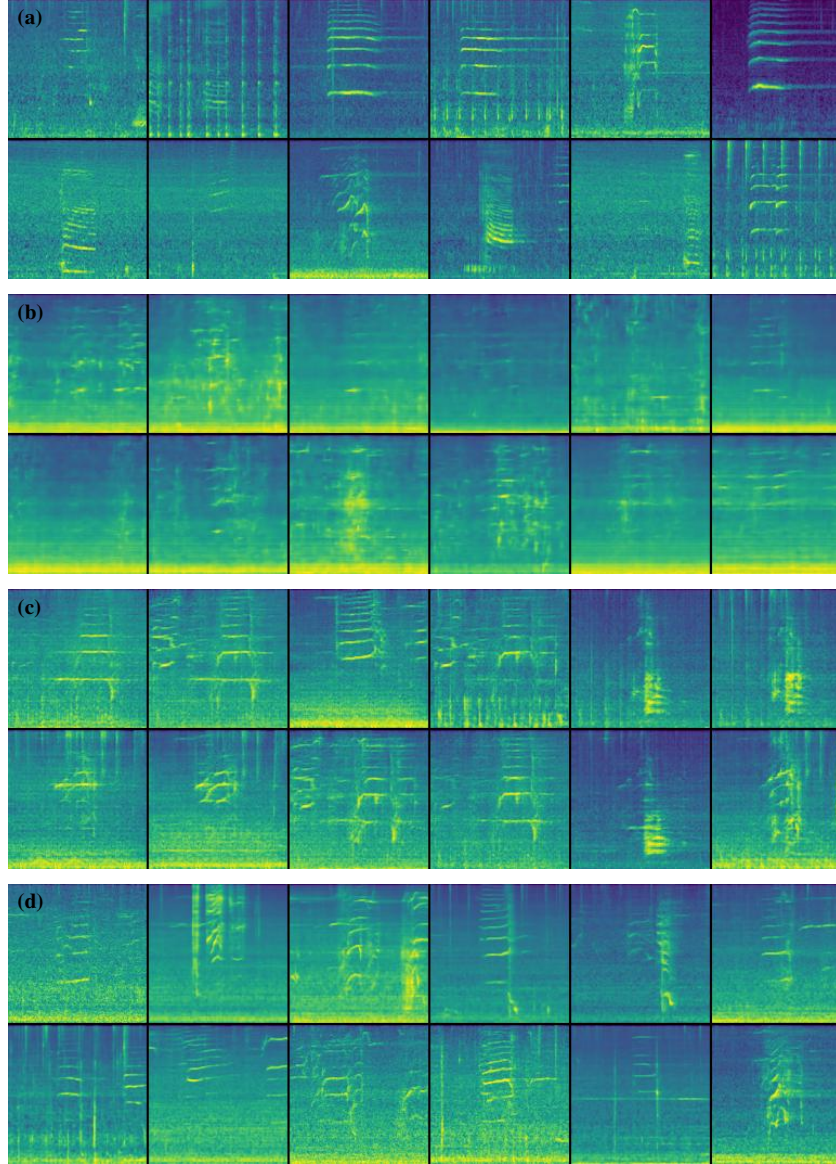


Figure 5: Examples of (a) real SRKW vocalizations, (b) VAE-generated samples, (c) GAN-generated samples, and (d) DDPM-generated samples.

independent test site.

#### 4. Discussion

The performance of deep learning models in bioacoustic classification tasks is strongly dependent on the availability and diversity of annotated training data. To address this limitation without relying on additional expensive and time consuming manual labeling, data augmentation is often employed to expand the effective size and variability of training sets. In this work, traditional data augmentation techniques were first evaluated for their ability to enhance model performance. The baseline model, trained exclusively on real SRKW vocalizations from the Lime Kiln dataset, performed poorly when tested on the independent Robert’s Bank site. This outcome was expected given the acoustic variability between locations and differences

in hydrophones. Differences in background noise profiles and sound propagation conditions often limit a model’s ability to generalize beyond its training domain. In contrast, augmenting the training data with time-shifted (experiment II) versions of the original vocalizations led to a substantial performance boost over the baseline model, highlighting its utility in scenarios where annotated vocalization samples are scarce. Time-shifting is computationally inexpensive, easy to implement, and resulted in consistently higher precision–recall performance compared to the baseline. Other lightweight augmentation methods such as pitch shifting (White et al., 2022; Duc, 2020), additive noise (White et al., 2022), and random masking (Park et al., 2019), have shown similarly measurable improvements across various bioacoustic applications.

More sophisticated approaches, such as vocalization mask augmentation (experiment III), further improved model perfor-



Table 2: Precision, Recall, and F1-Score at Threshold 0.5 for each augmentation strategy. Experiment I serves as the baseline with 1,261 samples. Other experiments include augmented data with increasing sample sizes. Values are reported as mean  $\pm$  standard deviation across ten runs. Bold values indicate the best performance within each augmentation level.

Experiment	Number of Samples	Precision	Recall	F1-Score
I ( <i>Baseline</i> )	1,261	0.65 $\pm$ 0.024	0.42 $\pm$ 0.048	0.51 $\pm$ 0.034
II ( <i>Time-shift</i> )	+ 5,000	0.70 $\pm$ 0.021	0.62 $\pm$ 0.030	0.66 $\pm$ 0.017
	+ 20,000	0.76 $\pm$ 0.020	0.71 $\pm$ 0.021	0.73 $\pm$ 0.012
III ( <i>Masks</i> )	+ 5,000	0.98 $\pm$ 0.004	0.51 $\pm$ 0.250	0.67 $\pm$ 0.021
	+ 20,000	<b>0.99</b> $\pm$ 0.003	0.53 $\pm$ 0.030	0.69 $\pm$ 0.025
IV ( <i>VAE</i> )	+ 5,000	0.69 $\pm$ 0.033	0.41 $\pm$ 0.057	0.51 $\pm$ 0.039
	+ 20,000	0.78 $\pm$ 0.039	0.45 $\pm$ 0.024	0.57 $\pm$ 0.019
V ( <i>GAN</i> )	+ 5,000	0.71 $\pm$ 0.025	0.41 $\pm$ 0.037	0.52 $\pm$ 0.026
	+ 20,000	0.76 $\pm$ 0.032	0.50 $\pm$ 0.051	0.60 $\pm$ 0.034
VI ( <i>DDPM</i> )	+ 5,000	0.70 $\pm$ 0.012	0.71 $\pm$ 0.040	0.71 $\pm$ 0.019
	+ 20,000	0.67 $\pm$ 0.029	<b>0.87</b> $\pm$ 0.024	0.75 $\pm$ 0.010
VII ( <i>Hybrid</i> )	+ 5,000	0.96 $\pm$ 0.011	0.63 $\pm$ 0.033	0.76 $\pm$ 0.021
	+ 20,000	<b>0.99</b> $\pm$ 0.003	0.69 $\pm$ 0.024	<b>0.81</b> $\pm$ 0.015

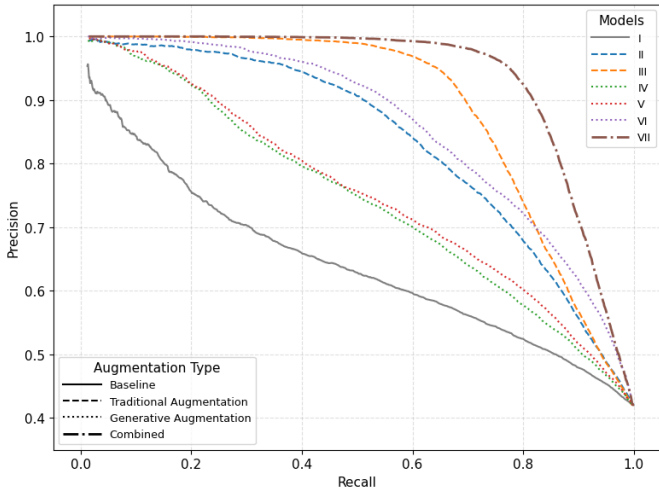


Figure 6: Mean Precision-Recall curves for all seven training regimes under the 20,000 augmentation setting, evaluated on the Robert’s Bank test set across the entire range of thresholds. Each curve shows the average performance over ten independent runs. The gray solid line represents the baseline model (I). Models II (blue dashed) and III (orange dashed) represent time-shifting and vocalization masking augmentations. Models IV (green dotted) and V (red dotted) and VI (purple dotted) use generative approaches based on VAE, GANs and DDPM, respectively. Finally, Model VII (brown dash-dot) combines the time-shift, masks and DDPM strategies.

mance by enabling context-aware signal insertion. By overlaying curated, high-quality vocalization templates onto real background noise from the target environment this method generated realistic and acoustically plausible training samples. Conceptually, this approach relates to prior work by Li et al. (2020), who synthesized dolphin whistle contours for model training by blending artificial whistle shapes with natural background spectrograms. However, our implementation differs by using empirically recorded SRKW vocalizations as masks, rather than synthetic contours, yielding ecologically grounded composites that preserve the spectral characteristics of genuine calls. Crucially, it exposed the model to the specific environmental context in which detection is expected to occur, allowing it to learn discriminative features that are directly relevant for deployment and can be observed by the drastic increase in precision, demonstrating the model’s improved ability to avoid false positives. However, because the curated call catalogue contained a limited number of vocalization types, and only high quality vocalizations, the model was not exposed to the full natural variability of SRKW vocalizations. In real-world conditions, such as those present in the Lime Kiln and Robert’s Bank datasets, vocalizations can often appear faint, overlapped, or embedded within complex acoustic environments. This is reflected by the more modest recall performance observed, indicating that the model failed to detect many calls due to its narrow exposure to vocalization diversity. Furthermore, increasing the number of augmented samples from +5,000 to +20,000 provided limited performance gains, suggesting that adding more homogeneous high-quality calls did not increase data variability or improve generalization.

Nonetheless, traditional augmentation techniques are ultimately limited by the scope and structure of the original dataset. There is only so much variability that can be meaningfully introduced through manual transformations. In this regard, gener-



ative approaches offer an alternative. Generative methods used in this study hold the potential to unlock new regions of the data distribution by synthesizing entirely novel yet realistic samples. All generative methods outperformed the baseline, indicating their ability to synthesize relevant and in-distribution samples that enhanced classifier performance. While both the VAE-based and GAN-based approaches underperformed traditional augmentation methods, they exhibited a more gradual decline in precision as recall increased. This pattern suggests that, despite some imperfections, generative models captured a broader spectrum of the vocalization space, including samples that were more challenging or ambiguous.

This trend is most evident with the DDPM-based approach, which outperformed the other generative models and time-shifting across the entire recall range, and surpassed all augmentation strategies at the upper end of the recall spectrum. More importantly, combining DDPM-generated samples with traditional augmentations led to the highest overall classification performance, exceeding what any single strategy could achieve in isolation. This hybrid configuration appears to be a promising strategy for bioacoustic applications, as it leverages the complementary strengths of generative and non-generative augmentations. The improvement is particularly notable when compared to the mask-only augmentation experiment where both the hybrid and mask-only models maintained similarly high precision at low recall levels, but only the hybrid approach sustained that precision as recall increased. These results suggest that generative and non-generative augmentations can complement each other. The former expands the diversity of vocalization forms beyond the limits of the original dataset, while the latter ensures contextual realism. Their combination enables the model to generalize more effectively by learning from a richer set of training scenarios.

Despite the promising results, several limitations must be acknowledged. Most notably, the use of generative models such as GANs, VAE, and DDPM in particular, introduces significant computational overhead. Training these models requires substantial GPU resources and time, particularly when compared to traditional augmentation techniques like time-shifting or template overlay, which are fast, lightweight, and easy to implement. Inference with DDPMs can also be relatively slow, especially for generating high quality samples, further limiting their practicality in real-time or large-scale training workflows. Additionally, contrary to CNN-based classifiers, which are now well established and supported by numerous off-the-shelf implementations, generative methods demand a certain degree of domain knowledge to implement and apply effectively. This higher complexity could limit their adoption, particularly among practitioners without specialized expertise. Moreover, while DDPM-based augmentation improved overall performance when used in combination with other methods, it did not consistently outperform conventional techniques when used in isolation. These findings broadly correspond with the work of [Herbst et al. \(2024\)](#), who reported that VAE and DDPM-based augmentation did not substantially improve classifier performance of Hainan gibbon calls beyond what could be achieved with a suite of traditional augmentation methods. This

suggests that the added complexity may not always be justified in scenarios where computational resources are constrained or rapid deployment is required. Thus, there is a clear tradeoff between the potential gains in data diversity offered by generative models and the practical demands of training and integrating them into existing pipelines.

Nonetheless, the ability of generative models to produce a potentially unlimited number of diverse samples, or at least far beyond what simpler augmentation methods can offer, holds significant promise for bioacoustic applications. The capacity to synthesize novel vocalization patterns, drawn from a real distribution, that still conform to the underlying structure of real vocalizations can complement training pipelines for models that need to generalize across ecotypes, acoustic environments, or rare call types.

## 5. Conclusion

In this work, we explored several data augmentation strategies to address the limitations of small, site-specific bioacoustic datasets and to improve generalization to new recording environments. Our results demonstrate that while traditional augmentation techniques such as time-shifting and vocalization masking offer tangible performance improvements with minimal computational overhead, generative models can further enrich training datasets and enhance model robustness. In contexts where expert annotation is difficult to obtain, or where vocalizations are rare due to limited accessibility (e.g., far offshore) or the scarcity of the species itself, these computational methods represent a justifiable investment, maximizing the utility of available data and improving overall model performance.

This study serves as a stepping stone toward the broader integration of generative models into bioacoustic data augmentation pipelines. While these models introduce added complexity and computational demands, our findings suggest that their inclusion in targeted augmentation strategies can be a worthwhile investment in scenarios where generalization and recall performance are critical. In particular, the combination of generative and traditional augmentation techniques appears to be a promising strategy for enhancing model robustness and cross-site generalization.

Future work could explore diffusion models that operate in a compressed latent space, where spectrograms are first encoded into a lower-dimensional representation before the generative process. This approach can dramatically reduce training and inference costs while preserving sample quality. Additionally, further research into quality control strategies, such as integrating filtering approach similarly to the a PCA-based strategy used in this study directly into the DDPM training process, could improve the reliability of generated samples and enhance the scalability and effectiveness of generative augmentation approaches.

## Funding sources

This research was supported by the Humans and Algorithms Listening and Looking for Orcas project (HALLO) funded by

the Canada Nature Fund for Aquatic Species at Risk of Fisheries and Oceans Canada (2022-NF-PAC-593022).

## Acknowledgments

We would like to thank April Houweling, Lauren Laturnus and Dr Jennifer Wladichuck for their help in annotating Killer Whale vocalizations.

## References

- Abeßer, J., 2020. A review of deep learning based methods for acoustic scene classification. *Applied Sciences* 10, 2020. doi:<https://doi.org/10.3390/app10062020>.
- Agarwal, S., Farid, H., 2021. Detecting deep-fake videos from aural and oral dynamics, in: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 981–989.
- Arjovsky, M., Bottou, L., 2017. Towards principled methods for training generative adversarial networks. *arXiv preprint arXiv:1701.04862*.
- Arjovsky, M., Chintala, S., Bottou, L., 2017. Wasserstein generative adversarial networks, in: *International conference on machine learning*, PMLR. pp. 214–223.
- Au, W.W., Ford, J.K., Horne, J.K., Allman, K.A.N., 2004. Echolocation signals of free-ranging killer whales (*orcinus orca*) and modeling of foraging for chinook salmon (*oncorhynchus tshawytscha*). *The Journal of the Acoustical Society of America* 115, 901–909.
- Baird, R.W., Stacey, P.J., 1988. Variation in saddle patch pigmentation in populations of killer whales (*orcinus orca*) from british columbia, alaska, and washington state. *Canadian Journal of Zoology* 66, 2582–2585. doi:<https://doi.org/10.1139/z88-380>.
- Barrett-Lennard, L.G., Ford, J.K., Heise, K.A., 1996. The mixed blessing of echolocation: differences in sonar use by fish-eating and mammal-eating killer whales. *Animal behaviour* 51, 553–565.
- Bergler, C., Barnhill, A., Perrin, D., Schmitt, M., Maier, A.K., Nöth, E., 2022. Orca-whisper: An automatic killer whale sound type generation toolkit using deep learning., in: *INTERSPEECH*, pp. 2413–2417.
- Bergler, C., Gebhard, A., Towers, J.R., Butyrev, L., Sutton, G.J., Shaw, T.J., Maier, A., Nöth, E., 2021. Fin-print a fully-automated multi-stage deep-learning-based framework for the individual recognition of killer whales. *Scientific reports* 11, 23480. doi:<https://doi.org/10.1038/s41598-023-38132-7>.
- Bergler, C., Schröter, H., Cheng, R.X., Barth, V., Weber, M., Nöth, E., Hofer, H., Maier, A., 2019. Orca-spot: An automatic killer whale sound detection toolkit using deep learning. *Scientific reports* 9, 10997. doi:<https://doi.org/10.1038/s41598-019-47335-w>.
- Binder, C.M., 2018. Impacts of environment-dependent acoustic propagation on passive acoustic monitoring of cetaceans. Ph.D. thesis. Dalhousie University.
- Brown, J.C., Smaragdis, P., Nousek-McGregor, A., 2010. Automatic identification of individual killer whales. *The Journal of the Acoustical Society of America* 128, EL93–EL98. doi:<https://doi.org/10.1121/1.3462232>.
- Cao, H., Tan, C., Gao, Z., Xu, Y., Chen, G., Heng, P.A., Li, S.Z., 2024. A survey on generative diffusion models. *IEEE Transactions on Knowledge and Data Engineering*.
- Che, T., Li, Y., Jacob, A.P., Bengio, Y., Li, W., 2016. Mode regularized generative adversarial networks. *arXiv preprint arXiv:1612.02136*.
- Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J., 2018. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8789–8797.
- Dede, A., Öztürk, A.A., Akamatsu, T., Tonay, A.M., Öztürk, B., 2014. Long-term passive acoustic monitoring revealed seasonal and diel patterns of cetacean presence in the istanbul strait. *Journal of the Marine Biological Association of the United Kingdom* 94, 1195–1202. doi:<https://doi.org/10.1017/S0025315413000568>.
- Dhariwal, P., Nichol, A., 2021a. Diffusion models beat gans on image synthesis, in: *Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., Vaughan, J.W. (Eds.), Advances in Neural Information Processing Systems*, Curran Associates, Inc. pp. 8780–8794. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/49ad23d1ec9fa4bd8d77d02681df5cfa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/49ad23d1ec9fa4bd8d77d02681df5cfa-Paper.pdf).
- Dhariwal, P., Nichol, A., 2021b. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* 34, 8780–8794.
- Donahue, C., McAuley, J., Puckette, M., 2018. Adversarial audio synthesis. *arXiv preprint arXiv:1802.04208*.
- Duc, P.N.H., 2020. Development of artificial intelligence methods for marine mammal detection and classification of underwater sounds in a weak supervision (but) Big Data-Expert context. Ph.D. thesis. Sorbonne Université.
- Elbir, A., Aydin, N., 2020. Music genre classification and music recommendation by using deep learning. *Electronics Letters* 56, 627–629. doi:<https://doi.org/10.1049/el.2019.4202>.
- Filatova, O., Fedutin, I., Nagaylik, M., Burdin, A., Hoyt, E., 2009. Usage of monophonic and biphonic calls by free-ranging resident killer whales (*orcinus orca*) in kamchatka, russian far east. *Acta ethologica* 12, 37–44. doi:<https://doi.org/10.1007/s10211-009-0056-7>.
- Ford, J.K., 1989. Acoustic behaviour of resident killer whales (*orcinus orca*) off vancouver island, british columbia. *Canadian Journal of Zoology* 67, 727–745.
- Ford, J.K., Ellis, G.M., Barrett-Lennard, L.G., Morton, A.B., Palm, R.S., Balcomb III, K.C., 1998. Dietary specialization in two sympatric populations of killer whales (*orcinus orca*) in coastal british columbia and adjacent waters. *Canadian journal of zoology* 76, 1456–1471. doi:<https://doi.org/10.1139/z98-089>.
- Ford, J.K., Fisher, H.D., 1978. Underwater acoustic signals of the narwhal (*monodon monoceros*). *Canadian Journal of Zoology* 56, 552–560. doi:<https://doi.org/10.1139/z78-079>.
- Ford, J.K., et al., 1987. A catalogue of underwater calls produced by killer whales (*orcinus orca*) in british columbia. URL: [https://www.researchgate.net/publication/285709635\\_A\\_catalogue\\_of\\_underwater\\_calls\\_produced\\_by\\_killer\\_whales\\_Orcinus\\_orca\\_in\\_British\\_Columbia](https://www.researchgate.net/publication/285709635_A_catalogue_of_underwater_calls_produced_by_killer_whales_Orcinus_orca_in_British_Columbia).
- Frazao, F., Joy, R., Dowd, M., 2025. Comparing acoustic representations for deep learning-based classification of underwater acoustic signals: A case study on orca (*orcinus orca*) vocalizations. *Ecological Informatics* 90, 103297. URL: <https://www.sciencedirect.com/science/article/pii/S1574954125003061>, doi:<https://doi.org/10.1016/j.ecoinf.2025.103297>.
- Gillespie, D., Caillat, M., Gordon, J., White, P., 2013. Automatic detection and classification of odontocete whistles. *The Journal of the Acoustical Society of America* 134, 2427–2437. doi:<https://doi.org/10.1121/1.4816555>.
- Gillespie, D., Mellinger, D., Gordon, J., McLaren, D., Redmond, P., McHugh, R., Trinder, P., Deng, X., Thode, A., 2008. Panguard: Semiautomated, open source software for real-time acoustic detection and localisation of cetaceans. *Journal of the Acoustical Society of America* 30, 54–62.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets, in: *Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K. (Eds.), Advances in Neural Information Processing Systems*, Curran Associates, Inc. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf).
- Gudivada, V., Apon, A., Ding, J., 2017. Data quality considerations for big data and machine learning: Going beyond data cleaning and transformations. *International Journal on Advances in Software* 10, 1–20.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C., 2017. Improved training of wasserstein gans. *Advances in neural information processing systems* 30.
- Hauer, C., Nöth, E., Barnhill, A., Maier, A., Guthunz, J., Hofer, H., Cheng, R.X., Barth, V., Bergler, C., 2023. Orca-spy enables killer whale sound source simulation, detection, classification and localization using an integrated deep learning-based segmentation. *Scientific Reports* 13, 11106. doi:<https://doi.org/10.1038/s41598-023-38132-7>.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- Herbst, C., Jeantet, L., Dufourq, E., 2024. Empirical evaluation of variational

- autoencoders and denoising diffusion models for data augmentation in bioacoustics classification, in: Annual Conference of South African Institute of Computer Scientists and Information Technologists, Springer. pp. 45–61. doi:[https://doi.org/10.1007/978-3-031-64881-6\\_3](https://doi.org/10.1007/978-3-031-64881-6_3).
- Ho, J., Jain, A., Abbeel, P., 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems* 33, 6840–6851.
- Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A., 2017. Image-to-image translation with conditional adversarial networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1125–1134.
- Kahl, S., Wood, C.M., Eibl, M., Klinck, H., 2021. Birdnet: A deep learning solution for avian diversity monitoring. *Ecological Informatics* 61, 101236. doi:<https://doi.org/10.1016/j.ecoinf.2021.101236>.
- Karaoglu, S., Gevers, T., et al., 2021. Self-supervised face image manipulation by conditioning gan on face decomposition. *IEEE Transactions on Multimedia* 24, 377–385.
- Karras, T., Laine, S., Aila, T., 2019. A style-based generator architecture for generative adversarial networks, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 4401–4410.
- King, S.L., Harley, H.E., Janik, V.M., 2014. The role of signature whistle matching in bottlenose dolphins, *tursiops truncatus*. *Animal Behaviour* 96, 79–86. doi:<https://doi.org/10.1016/j.anbehav.2014.07.019>.
- Kingma, D.P., 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* doi:<https://doi.org/10.48550/arXiv.1312.6114>.
- Kingma, D.P., Welling, M., et al., 2019. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning* 12, 307–392.
- Kirsebom, O.S., Frazao, F., Simard, Y., Roy, N., Matwin, S., Giard, S., 2020. Performance of a deep neural network at detecting north atlantic right whale upcalls. *The Journal of the Acoustical Society of America* 147, 2636–2646. doi:[10.1121/10.0001132](https://doi.org/10.1121/10.0001132).
- Kong, Z., Ping, W., Huang, J., Zhao, K., Catanzaro, B., 2020. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *nature* 521, 436–444. doi:<https://doi.org/10.1038/nature14539>.
- Li, L., Qiao, G., Liu, S., Qing, X., Zhang, H., Mazhar, S., Niu, F., 2021. Automated classification of tursiops aduncus whistles based on a depth-wise separable convolutional neural network and data augmentation. *The Journal of the Acoustical Society of America* 150, 3861–3873. doi:<https://doi.org/10.1121/10.0007291>.
- Li, P., Liu, X., Palmer, K., Fleishman, E., Gillespie, D., Nosal, E.M., Shiu, Y., Klinck, H., Cholewiak, D., Helble, T., et al., 2020. Learning deep models from synthetic data for extracting dolphin whistle contours, in: 2020 International Joint Conference on Neural Networks (IJCNN), IEEE. pp. 1–10. doi:[10.1109/IJCNN48605.2020.9206992](https://doi.org/10.1109/IJCNN48605.2020.9206992).
- Li, P., Roch, M.A., Klinck, H., Fleishman, E., Gillespie, D., Nosal, E.M., Shiu, Y., Liu, X., 2023. Learning stage-wise gans for whistle extraction in time-frequency spectrograms. *IEEE Transactions on Multimedia* 25, 9302–9314. doi:[10.1109/TMM.2023.3251109](https://doi.org/10.1109/TMM.2023.3251109).
- Manilow, E., Seetharman, P., Salamon, J., 2020. Open Source Tools & Data for Music Source Separation. <https://source-separation.github.io/tutorial>. URL: <https://source-separation.github.io/tutorial>.
- Mariani, G., Scheidegger, F., Istrate, R., Bekas, C., Malossi, C., 2018. Bagan: Data augmentation with balancing gan. *arXiv preprint arXiv:1803.09655*.
- Miller, P.J., 2006. Diversity in sound pressure levels and estimated active space of resident killer whale vocalizations. *Journal of Comparative Physiology A* 192, 449–459. doi:<https://doi.org/10.1007/s00359-005-0085-2>.
- Mishachandar, B., Vairamuthu, S., 2021. Diverse ocean noise classification using deep learning. *Applied Acoustics* 181, 108141. doi:<https://doi.org/10.1016/j.apacoust.2021.108141>.
- Morfi, V., Lachlan, R.F., Stowell, D., 2021. Deep perceptual embeddings for unlabelled animal sound events. *The Journal of the Acoustical Society of America* 150, 2–11. doi:<https://doi.org/10.1121/10.0005475>.
- Morin, P.A., McCarthy, M.L., Fung, C.W., Durban, J.W., Parsons, K.M., Perrin, W.F., Taylor, B.L., Jefferson, T.A., Archer, F.I., 2024. Revised taxonomy of eastern north pacific killer whales (*orcinus orca*): Bigg’s and resident ecotypes deserve species status. *Royal Society Open Science* 11, 231368. doi:<https://doi.org/10.1098/rsos.231368>.
- Murphy, D.T., Ioup, E., Hoque, M.T., Abdelguerfi, M., 2022. Residual learning for marine mammal classification. *IEEE Access* 10, 118409–118418. doi:[10.1109/ACCESS.2022.3220735](https://doi.org/10.1109/ACCESS.2022.3220735).
- Nichol, A.Q., Dhariwal, P., 2021. Improved denoising diffusion probabilistic models, in: International conference on machine learning, PMLR. pp. 8162–8171.
- Nieto-Mora, D.A., Ferreira de Oliveira, M.C., Sanchez-Giraldo, C., Duque-Muñoz, L., Isaza-Narváez, C., Martínez-Vargas, J.D., 2024. Soundscape characterization using autoencoders and unsupervised learning. *Sensors* 24, 2597. doi:<https://doi.org/10.3390/s24082597>.
- Olson, J.K., Wood, J., Osborne, R.W., Barrett-Lennard, L., Larson, S., 2018. Sightings of southern resident killer whales in the salish sea 1976–2014: the importance of a long-term opportunistic dataset. *Endangered Species Research* 37, 105–118. doi:<https://doi.org/10.3354/esr00918>.
- Ozer, I., Ozer, Z., Findik, O., 2018. Noise robust sound event classification with convolutional neural network. *Neurocomputing* 272, 505–512. doi:<https://doi.org/10.1016/j.neucom.2017.07.021>.
- Padovese, B., Frazao, F., Kirsebom, O.S., Matwin, S., 2021. Data augmentation for the classification of north atlantic right whales upcalls. *The Journal of the Acoustical Society of America* 149, 2520–2530. doi:<https://doi.org/10.1121/10.0004258>.
- Padovese, B., Kirsebom, O.S., Frazao, F., Evers, C.H., Beslin, W.A., Theriault, J., Matwin, S., 2023. Adapting deep learning models to new acoustic environments—a case study on the north atlantic right whale upcall. *Ecological Informatics* 77, 102169. doi:<https://doi.org/10.1016/j.ecoinf.2023.102169>.
- Palmer, K., Cummings, E., Dowd, M.G., Frasier, K., Frazao, F., Harris, A., Houweling, A., Kanes, J., Kirsebom, O.S., Klinck, H., et al., 2025. A public dataset of annotated orcinus orca acoustic signals for detection and ecotype classification. *Scientific Data* 12, 1137.
- Park, D.S., Chan, W., Zhang, Y., Chiu, C.C., Zoph, B., Cubuk, E.D., Le, Q.V., 2019. SpecAugment: A simple data augmentation method for automatic speech recognition, in: Proc. Interspeech 2019, pp. 2613–2617. doi:[10.21437/Interspeech.2019-2680](https://doi.org/10.21437/Interspeech.2019-2680).
- Priestley, M., O’donnell, F., Simperl, E., 2023. A survey of data quality requirements that matter in ml development pipelines. *ACM Journal of Data and Information Quality* 15, 1–39. doi:<https://doi.org/10.1145/3592616>.
- Rabiner, L., Juang, B.H., 1993. Fundamentals of speech recognition. Prentice-Hall, Inc.
- Radford, A., Metz, L., Chintala, S., 2015. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.
- Reichert, M.S., Ronacher, B., 2015. Noise affects the shape of female preference functions for acoustic signals. *Evolution* 69, 381–394. doi:<https://doi.org/10.1111/evo.12592>.
- Riesch, R., Barrett-Lennard, L.G., Ellis, G.M., Ford, J.K., Deecke, V.B., 2012. Cultural traditions and the evolution of reproductive isolation: ecological speciation in killer whales? *Biological Journal of the Linnean Society* 106, 1–17. doi:<https://doi.org/10.1111/j.1095-8312.2012.01872.x>.
- Riesch, R., Ford, J.K., Thomsen, F., 2008. Whistle sequences in wild killer whales (*orcinus orca*). *The Journal of the Acoustical Society of America* 124, 1822–1829. doi:<https://doi.org/10.1121/1.2956467>.
- Roch, M.A., Miller, P., Helble, T.A., Baumann-Pickering, S., Širović, A., 2017. Organizing metadata from passive acoustic localizations of marine animals. *The Journal of the Acoustical Society of America* 141, 3605–3605. doi:<https://doi.org/10.1121/1.4987714>.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III 18, Springer. pp. 234–241.
- Sato, M., Trites, A.W., Gauthier, S., 2021. Southern resident killer whales encounter higher prey densities than northern resident killer whales during summer. *Canadian Journal of Fisheries and Aquatic Sciences* 78, 1732–1743. doi:<https://doi.org/10.1139/cjfas-2020-0445>.
- Shapiro, A.D., Wang, C., 2009. A versatile pitch tracking algorithm: From human speech to killer whale vocalizations. *The Journal of the Acoustical Society of America* 126, 451–459. doi:<https://doi.org/10.1121/1.3132525>.
- Sharpe, D.L., Castellote, M., Wade, P.R., Cornick, L.A., 2019. Call types of bigg’s killer whales (*orcinus orca*) in western alaska: Using vocal dialects to assess population structure. *Bioacoustics* 28, 74–99. doi:<https://doi.org/10.1080/09524622.2017.1396562>.
- Shim, J.Y., Kim, J., Kim, J.K., 2021. S2i-bird: Sound-to-image generation of bird species using generative adversarial networks, in: 2020 25th

- International Conference on Pattern Recognition (ICPR), pp. 2226–2232. doi:[10.1109/ICPR48806.2021.9412721](https://doi.org/10.1109/ICPR48806.2021.9412721).
- Shiu, Y., Palmer, K., Roch, M.A., Fleishman, E., Liu, X., Nosal, E.M., Helble, T., Cholewiak, D., Gillespie, D., Klinck, H., 2020. Deep neural networks for automated detection of marine mammal species. *Scientific reports* 10, 607. doi:<https://doi.org/10.1038/s41598-020-57549-y>.
- Shorten, C., Khoshgoftaar, T.M., 2019. A survey on image data augmentation for deep learning. *Journal of big data* 6, 1–48. doi:<https://doi.org/10.1186/s40537-019-0197-0>.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S., 2015. Deep unsupervised learning using nonequilibrium thermodynamics, in: *International conference on machine learning*, pmlr. pp. 2256–2265.
- Srivastava, A., Valkov, L., Russell, C., Gutmann, M.U., Sutton, C., 2017. Vee-gan: Reducing mode collapse in gans using implicit variational learning. *Advances in neural information processing systems* 30.
- Stowell, D., 2022. Computational bioacoustics with deep learning: a review and roadmap. *PeerJ* 10, e13152. doi:<https://doi.org/10.7717/peerj.13152>.
- Thomsen, F., Franck, D., Ford, J., 2001. Characteristics of whistles from the acoustic repertoire of resident killer whales (*Orcinus orca*) off vancouver island, british columbia. *The Journal of the Acoustical Society of America* 109, 1240–1246. doi:<https://doi.org/10.1121/1.1349537>.
- Tiwari, R.G., Gautam, V., Trivedi, N.K., Jain, A.K., Sharma, V., 2023. A deep learning approach for marine animal classification: Enhancing taxonomic identification and conservation efforts, in: *2023 3rd Asian Conference on Innovation in Technology (ASIANCON)*, pp. 1–8. doi:[10.1109/ASIANCON58793.2023.10270176](https://doi.org/10.1109/ASIANCON58793.2023.10270176).
- Trabucco, B., Doherty, K., Gurinas, M., Salakhutdinov, R., 2023. Effective data augmentation with diffusion models. *arXiv preprint arXiv:2302.07944* doi:<https://doi.org/10.48550/arXiv.2302.07944>.
- Wang, D., Chen, J., 2018. Supervised speech separation based on deep learning: An overview. *IEEE/ACM transactions on audio, speech, and language processing* 26, 1702–1726. doi:[10.1109/TASLP.2018.2842159](https://doi.org/10.1109/TASLP.2018.2842159).
- Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B., 2018. High-resolution image synthesis and semantic manipulation with conditional gans, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8798–8807.
- Wang, X., Wang, K., Lian, S., 2020. A survey on face data augmentation for the training of deep neural networks. *Neural computing and applications* 32, 15503–15531.
- Webster, M.S., Budney, G.F., 2017. Sound archives and media specimens in the 21st century. *Comparative bioacoustics: An overview* , 479–503.
- White, E.L., White, P., Bull, J., Risch, D., Beck, S., Edwards, E., 2022. More than a whistle: Automated detection of marine sound sources with a convolutional neural network. *Frontiers in Marine Science* 9.
- Zhang, L., Huang, H.N., Yin, L., Li, B.Q., Wu, D., Liu, H.R., Li, X.F., Xie, Y.L., 2022. Dolphin vocal sound generation via deep wavegan. *Journal of Electronic Science and Technology* 20, 100171. doi:<https://doi.org/10.1016/j.jnlest.2022.100171>.