# Learning Stage-wise GANs for Whistle Extraction in Time-Frequency Spectrograms

Pu Li[1,2], Marie A. Roch[1], Holger Klinck[3,4], Erica Fleishman[5], Douglas Gillespie[6], Eva-Marie Nosal[7], Yu Shiu[8], and Xiaobai Liu[1], *Member, IEEE*

*Abstract*—Whistle contour extraction aims to derive animal whistles from time-frequency spectrograms as polylines. For toothed whales, whistle extraction results can serve as the basis for analyzing animal abundance, species identity, and social activities. During the last few decades, as long-term recording systems have become affordable, automated whistle extraction algorithms were proposed to process large volumes of recording data. Recently, a deep learning-based method demonstrated superior performance in extracting whistles under varying noise conditions. However, training such networks requires a large amount of labor-intensive annotation, which is not available for many species. To overcome this limitation, we present a framework of stage-wise generative adversarial networks (GANs), which compile new whistle data suitable for deep model training via three stages: generation of background noise in the spectrogram, generation of whistle contours, and generation of whistle signals. By separating the generation of different components in the samples, our framework composes visually promising whistle data and labels even when few expert annotated data are available. Regardless of the amount of human-annotated data, the proposed data augmentation framework leads to a consistent improvement in performance of the whistle extraction model, with a maximum increase of 1.69 in the whistle extraction mean F1-score. Our stage-wise GAN also surpasses one single GAN in improving whistle extraction models with augmented data. The data and code will be available at https://github.com/Paul-LiPu/CompositeGAN_WhistleAugment.

*Index Terms*—Data Augmentation, Generative Adversarial Networks.

## I. INTRODUCTION

### A. Background

Spectrograms in the time × frequency domain can show signal structure and are frequently used in audio analysis [1] . Patterns in spectrograms are used for sound event classification [2], bird song recognition [3], music genre classification [4], automatic music transcription [5], speech emotion recognition [6], and other tasks. Many acoustic signals have frequency-modulated (FM) components that are visible in spectrograms. Examples include human speech [7], human singing [8], cries of newborns [9], vocal melodies [10], and whale calls [11]. In this paper, we concentrate on whistles, the characteristic FM tonal calls of toothed whales.

Whale calls are used to study species identity [12] [13], individual identity [14], behavior [15] [16], communication and social activities [17], and density and abundance [18].

[1] San Diego State University [2] University of California, Irvine [3] Cornell University [4] Oregon State University [5] Colorado State University [6] University of St Andrews [7] University of Hawaii [8] RedRoute, Inc
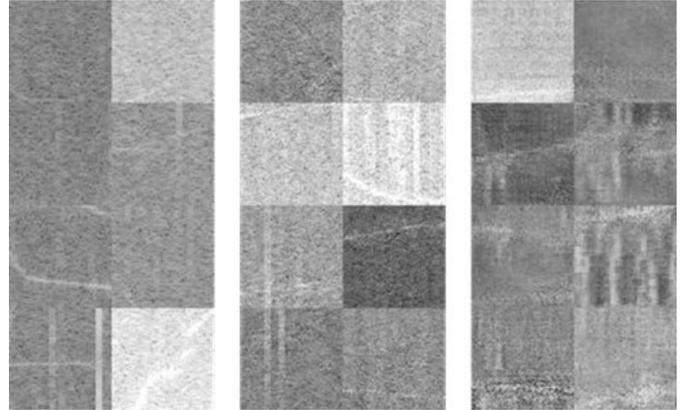


Fig. 1. Examples of spectrogram patches of (i) real samples (left); (ii) samples generated by our stage-wise GAN; (iii) samples generated by a single GAN. Multiple 64×64 patches are concatenated for better visualization.

Because whistles appear in spectrograms as characteristic contour shapes (Fig. 1 top left), experts can manually recognize animals' occurrence and label whistles as polylines on spectrograms. Whistle extraction algorithms [11] [19] [20] [21] [22] [23] aim to automate this process and identify each whistle as a polyline. Such extraction is challenging because of the high spatial and temporal variation of ocean sounds. Signals to be analyzed can be affected by recording device characteristics, sea state and propagation conditions, animal behavior, vocalizations from non-target species, and anthropogenic sounds, such as shipping and sonar.

Traditional methods (e.g., graph search [11]) first extract the spectral peaks, i.e., bins with local maximum energy on spectrogram, and then track the trace of whistle signals on the spectrogram by polynomial fitting of peaks [11] or probabilistic modeling [22]. Recently, [19] adapted convolutional neural networks (CNNs) to extract whistles and achieved improved performance. Instead of using spectral peaks, [19] predicts the confidence associated with the probability that a whistle signal appears in each time × frequency bin, which is similar to semantic segmentation in computer vision. [19] then uses graph search [11] to connect bins that are likely to contain a signal. By learning from a large set of annotated samples, the whistle extraction model can recognize noise and whistle patterns, and improve on graph search and probabilistic model results by a large margin. However, the performance of learning-based methods may degrade significantly as the amount of annotated data decreases, and large datasets are not always available because whistle annotation is expensive

and time-consuming. This motivates us to explore ways to synthesize whistle data cheaply with existing data by applying learning-based data generation methods.

### B. Objectives

The primary focus of this work was to develop methods that improve whistle extraction models when data are limited, thereby reducing the amount of data annotation required to recognize whistles. Therefore, our experiments mainly addressed situations with few data, and we sought to mitigate the effect of overfitting and improve the model's transferability for recognizing tonal signals. Although there are many ways to reduce overfitting, e.g., semi-supervised learning [24] and regularization [25], we focused on data augmentation methods for two reasons. First, we seek a method that can be applied to all datasets of frequency-modulated signal, including those containing no unannotated data. Semi-supervised learning may not be applicable in this scenario. Second, we are interested in characterizing the distribution of whistle data and exploring the effect of novel data on extraction of tonal signals. Regularization terms may not provide insight in this context. We note that our data augmentation method may be combined with a semi-supervised framework or loss function regularization to further improve the system performance. Though it is interesting to have these techniques involved, it is beyond the scope of this work.

Common audio or image data augmentation methods usually transform existing data to acquire new data, e.g., by adding Gaussian noise [26], and the augmented samples may implicitly act as a regularizer for the training of deep models [27]. But the distribution of the augmented samples may not be similar to that of the original data; e.g., generated whistle data may have abnormal contour shapes or unrealistic background noise. Previous work [19] generated novel samples by adding whistle contours to negative samples (background noise that contains no whistle signals), which simulated the situation where the same whistles occur in different ocean environments. However, the generated data did not include novel whistle shapes or background noise patterns, which restricted the variance in the data.

In this paper, our goal is to generate novel pairs of whistle data and labels. Although changes in noise affect vocalizations of many taxa [28], including toothed whales [29], we make the simplifying assumption that background noise is independent of whistle contours (contour-shape segmentation of whistles, which indicates the location of whistles on spectrogram and the whistles' frequency modulation). On the basis of this assumption, we decouple the synthesis of background noise and whistle contours. The generated whistle contours are used as labels for the model in [19]. Next, we add generated whistle signals with the desired contour shape to the spectrogram of background noise; i.e., we generate corresponding whistle data for the whistle contour.

We design our whistle generation algorithm as a series of three generative adversarial network (GAN) modules. The first GAN learns the ocean noise environment; it maps random numbers that have a Gaussian distribution to spectrograms representing background noise. The next GAN learns to map random inputs to spectrograms with whistle-like FM sweeps. The third GAN combines the outputs of the first two GAN modules, synthetic background noises and whistles, to obtain a synthetic whistle spectrogram. The generated whistle should follow the whistle contour's shape in the input. We employ an unpaired domain transfer framework, CycleGAN [30], to learn how synthetic noise and whistles can be merged into a synthetic spectrogram. While the original CycleGAN can generate slightly misaligned whistle signals from the desired contours, we exploit the whistle extraction network learned from annotated data to enforce the bin-wise consistency between generated whistles and input contours.

Another challenge is that GANs may not learn well with limited data. This may lead to corrupted synthesis, especially of the whistle contours. We observe that corrupted data have less confident predictions: the predicted probability is neither close to 0 nor close to 1, and thus the entropy is high. Accordingly, we introduce a method to prune such low-quality generated samples. Furthermore, because imperfect learning by GANs with few data may lead to discrepancies between the distributions of real data and generated data, we employ auxiliary batch normalization (ABN) layers [31] which separate the statistics of real and generated data to reduce the possible harmful effect of training with generated data.

### C. Contributions

We made three contributions. First, we proposed the stage-wise composite GANs to generate novel whistle extraction data, including spectrograms and corresponding whistle contour labels. Our experiments showed that the proposed stage-wise GAN surpassed the vanilla GANs with respect to the visual quality of the generated data (Fig. 1 middle and right). Second, we designed a comprehensive strategy to use GAN-generated samples to improve whistle-extraction models. We set criteria to remove corrupted data and we redesigned the whistle extraction network by adding ABN layers to optimize the training with generated data. Third, we applied our proposed data augmentation methods to varied amounts of whistle extraction data and observed consistent and significant improvements. Although GAN frameworks have been used for spectrogram generation and data augmentation in audio recognition tasks [32], to our knowledge, this is the first work to apply GAN-based augmentation to audio spectrogram segmentation data.

## II. RELATED WORKS

### A. Whistle Contour Extraction

There are three main classes of methods for extracting whale frequency-modulated whistles. The first is models that predict the probability of whistle peaks conditioned on past observations. Examples of this class include tests of hypothesized spectrogram region distributions [33], Bayesian inference [34], Kalman filters [20], and Monte-Carlo density filters [11] [35] [36] [22]. The second class, trajectory-search methods, seeks energy peaks along the frequency dimension and connects those peaks along the time dimension on the basis of trajectory
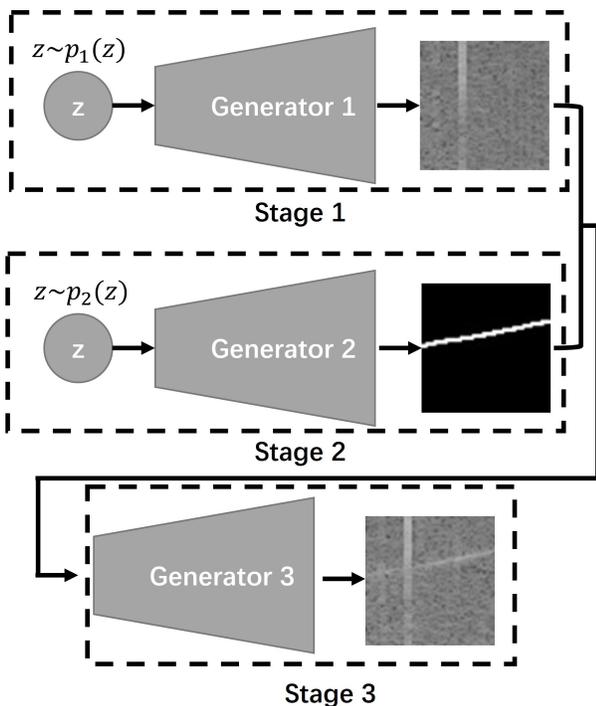
Fig. 2. Sketch of the proposed stage-wise GAN frameworks. The first two generators produce a spectrogram patch of background noise and a spectrogram patch of foreground whistle contour, respectively. These patches serve as inputs for the third generator.

estimation [12] [11] [37]. Improved trajectory-search methods reduce excessive numbers of false positives by applying ridge regression to local contexts [38] or energy minimization algorithms to ridge regression maps [39].

In recent years, the third class, deep learning methods, has been applied to process tonal information. Early works included extraction of information from human speech [40] and music [41]. Deep neural networks were also applied to toothed whale whistles [13] [42], but the goal of these works was to classify a time segment to species or call type rather than to extract detailed time × frequency information. In [19], we proposed a deep neural network to extract time × frequency contours of individual whistles. We apply our proposed data augmentation system to the training of whistle extraction model developed in [19].

### B. Generative adversarial networks

Generative adversarial networks (GANs) are a category of generative models. GANs are widely used for artificial image generation, e.g., face manipulation [43], compression noise removal [44], and generating images of people [45]. We adapted the methods from these computer vision tasks to generate realistic spectrograms that served as our novel training data. The landmark work on GANs [46] proposed one generator network that synthesizes samples ($G : X \rightarrow Y$, where $X$ is a random vector and $Y$ is a generated sample) and one discriminator network that learns to distinguish between generated samples and real samples. These networks are coupled in a zero-sum game with each network trying

to outperform the other. Following [46], researchers have improved the network architecture of GANs [47] and objective functions [48] to stabilize the training of GANs. Those GANs implicitly learn the distribution of real samples, and novel data can be sampled from the distribution. We employ this type of GAN to generate novel spectrogram noises and whistle contours.

Another type of GAN tackles the image-to-image translation problem, aiming to learn a mapping ($F : x \rightarrow y$, where $x \in X$, $y \in Y$) between a source domain $X$ and a target domain $Y$, e.g., transfer a horse in the image to a zebra. CycleGAN [30] extends this idea by leaning two mappings ($F : x \rightarrow y$, $G : y \rightarrow x$, where $x \in X$, $y \in Y$) without the need for pairwise correspondence between the elements of $X$ and $Y$. This idea can be adapted to our task to generate spectrograms containing whistles, where $X$ is the domain containing pairs of desired whistle contours and spectrograms with background noise, and $Y$ consists of spectrograms with whistles and noise. Recent work improved the idea of [30] by adding a spatial attention mechanism [49] and image quality assessment term [50].

### C. GAN-based augmentation

GANs provide an option to generate novel data by learning the distribution of existing data and sampling data from the distribution, which is a valuable addition to the common augmentation techniques that are based on data transformation. Vanilla GAN models, which map random numbers to generated samples, have been used for data augmentation. [51] trained a GAN model to augment computed tomography (CT) images of livers for the classification of lesions. [52] applied a conditional GAN to augment samples from given categories and restore the balance of imbalanced image classification data. [53] applied progressively grown GANs (PGGANs) to a brain segmentation task, and the generator learned to synthesize the generated sample and corresponding segmentation labels. Domain transfer GANs have also been used for data augmentation. [54] applied CycleGAN to day-to-night image translation, which helped to improve the object detection model.

Despite the success of GANs in synthesizing visually appealing samples and augmenting existing data, there are still limitations of GANs for synthesizing high-quality augmented data, especially for pixel-wise regression tasks such as semantic segmentation. First, GANs usually suffer from mode collapses [55]: the generated samples may have lower variance than the real samples. Second, GANs may generate samples with artifacts or failure regions [56], which may especially hamper the training of pixel-wise regression tasks. A sample selection method may be required to choose high-quality samples from the GAN-generated samples [57]. Third, the training of GANs can be unstable, which results in different distributions of generated samples and real samples [58].

Therefore, GAN-based data augmentation usually requires improving the quality of generated samples. A common solution is to use real samples or computer graphics models in the generator network. In [59], the GAN learned to generate

samples conditioned on real samples and random numbers. Similarly, [60] transferred synthetic images built by computer graphics models to realistic images, and the augmented samples improved models in estimation of gazes, hand poses, and animal poses. Another way to improve training of the GAN is to use supervision from target tasks. [61] added an auxiliary classification head on the discriminator of GAN and used the classification loss to guide discriminator and generator learning.

Recently, stage-wise GANs were proposed to augment data for pixel-wise regression tasks. [62] employed a two-stage GAN augmentation of cell nuclei segmentation data. Their framework generates a cell nuclei segmentation mask in the first stage and images of nuclei in the second stage. Our proposed method is closely related to [62], and we further separate the learning of object appearance and the segmentation mask. This separation can be extended to other semantic segmentation scenarios. For example, when generating a scene containing road and cars, our framework may first generate the appearance of the road and car independent of the segmentation mask, then generate an image of the scene according to the segmentation mask and the appearance of the objects (road and cars). In this way, our framework explores the distribution of object appearance and provides variance in the appearance of objects in the generated image of the scene. Another improvement is that we employ the knowledge from segmentation networks to regularize bin-wise correspondence between generated samples and labels.

## III. METHODS

The objective of this work is to develop a data augmentation approach to generate novel data for whistle extraction. We treat the cropped patches from the time-frequency spectrograms as data samples, and we employ stage-wise GANs, which we call WAS-GANs (**W**histle **A**ugmentation **S**tage-wise **G**enerative **A**dversarial **N**etworks), to generate both negative samples (noise only) and positive samples (whistles in the presence of noise). Our techniques can be extended to other acoustic tasks or computer vision tasks, e.g., sound classification and semantic segmentation.

Fig. 2 illustrates the three stages of our sample generation approach. In Stage 1, a Wasserstein GAN with gradient penalty (WGAN-gp) [48] learns to produce the negative samples containing background noises. In Stage 2, we train another WGAN-gp model with the real whistle contour annotations to generate whistle contour segmentation masks. In Stage 3, we use a CycleGAN [30] to generate positive samples. The whistle signals are added to the negative samples obtained in Stage 1 according to contour shapes defined in Stage 2. The positive samples and segmentation masks are used as the whistle extraction data and labels, respectively. Both generated negative samples and positive samples are used to train the whistle extraction model, and the resulted whistle extraction performance is used to assess our GAN-based augmentation.

### A. GAN-based negative sample synthesis

We assume that the underwater background noise (negative samples) follows an implicit distribution. The generator learns
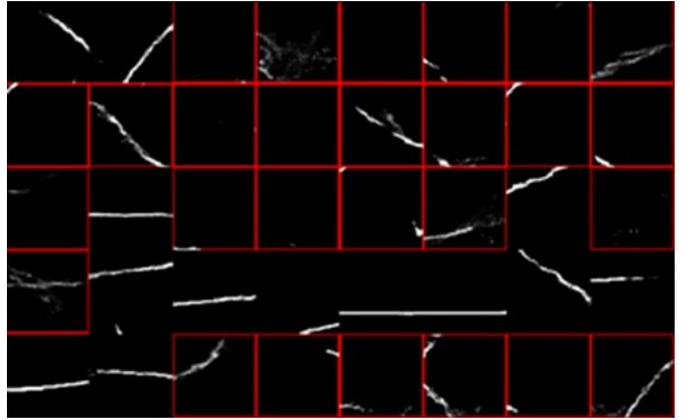


Fig. 3. Illustration of whistle contour selection. Low-quality generated patches are highlighted by red bounding boxes. Multiple $64\times64$ patches are concatenated.

the mapping between a multivariate Gaussian distribution and the distribution of negative samples. While many GAN models can learn this mapping, we chose WGAN-gp because its training is relatively stable [48]. The model includes a generator network, $G$, and a discriminator network, $D$. Network $G$ maps a multivariate Gaussian random variable to generate negative samples. Network $D$ estimates the Wasserstein distance between real samples and generated background noise (negative) samples. We denote $P_r$ as the distribution of real data $x$; $P_g$ as the distribution of generated data implicitly defined by $\widetilde{x} = G(z)$, where $z$ is a random vector following the standard multivariate Gaussian distribution; and $\hat{x}$ as a randomly weighted sum of x and $\widetilde{x}$. The loss function for the discriminator network is defined as

$$L = \mathbb{E}_{\tilde{x}\sim\mathbb{P}_g}\ [D\,(\tilde{x})] - \mathbb{E}_{x\sim\mathbb{P}_r}\ [D\,(x)] + \\ \lambda\mathbb{E}_{\hat{x}\sim\mathbb{P}_{\hat{x}}}\ \left[(||\,\nabla_{\hat{x}}D\,(\hat{x})\,||_2 - 1)^2\right] \quad (1)$$

where $\nabla_{\hat{x}}D\,(\hat{x})$ is the gradient of discriminator $D$'s output on $\hat{x}$. This loss function encourages the discriminator to maximize the estimated Wasserstein distance between real and generated samples. The gradient penalty term $\mathbb{E}_{\hat{x}\sim\mathbb{P}_{\hat{x}}}\ \left[(||\,\nabla_{\hat{x}}D\,(\hat{x})\,||_2 - 1)^2\right]$ enforces a soft version of Lipschitz constraint on the discriminator network. The loss function for the generator network is

$$L_G = \mathbb{E}_z\ [-D\,(G(z))] \quad (2)$$

which encourages the generator to generate samples that have a small estimated Wasserstein distance from the real samples, i.e., to follow a distribution similar to that of the real data.

### B. GAN-based positive sample synthesis

We split synthesis of positive samples (spectrograms containing whistles) into two stages: generation of whistle contours and injection of the whistle into synthetic background noise. In the first stage, we employ the same networks and loss functions as in Section III-A, given the assumption that the shape of whistle contours is independent of the underwater environments.

In the second stage, we aim to generate positive samples according to the synthetic background noise and whistle contours. We treat this as an unpaired domain transfer task, which can be solved effectively by CycleGAN [30]. Our source domain, A, consists of pairs of negative samples and whistle contours, and the target domain B includes positive samples. We adopt the CycleGAN from [30] for our experiments, but any improved model readily can be used in our framework.

There are two sets of generator and discriminator networks in CycleGAN. $G_A$ denotes the generator network that transfers samples from domain $A$ to domain $B$, i.e., generates whistle with the desired shape on the background noise spectrogram. $D_A$ denotes the discriminator network that distinguishes between real and generated spectrograms in domain $B$. $G_B$ is the network that transfers samples from domain $B$ to domain $A$, effectively separating the whistle contour from the background noise. Because we assume that the whistle contour and background noise are independent, we do not use a single $D_B$ network for the joint distribution of whistle contours and background noise. Instead, we use two $D_B$ networks for the marginal distributions, one to discriminate negative samples and one to discriminate whistle contours.

Instead of directly generating positive samples by $G_A$, we let $G_A$ predict a residual term (whistle signals without background noises) to be added to the negative samples. By denoting a negative sample as $I_N$, a whistle contour as $I_W$, and the generated positive sample as $I'_P$, the process can be described as

$$I'_P = I_N + \gamma G_A(I_N, \ I_W) \tag{3}$$

where $\gamma$ is a factor that controls the signal strength and accounts for variability in the received signal level. This parameter can simulate the variation in signal strength caused by variation in signal source strength or the distance between the animal and recording devices.

To enforce the bin-wise correspondence between generated positive samples and whistle contours, i.e., to avoid misalignment between generated whistle extraction data and labels, we use the whistle extraction models, which are trained on the same set of real samples as CycleGAN, to design a regularization term for $G_A$ training. We call this term a loss function for the pixel-wise consistency, and represent it as

$$L_{consistence} = ||f(I'_P) - I_W||_1 \tag{4}$$

where $f$ denotes the whistle extraction model and $f(x)$ is the model's output, a confidence map indicating the presence of whistle energy in each bin of the spectrogram, with an input $x$. This loss encourages the whistle signals to appear at the same position as the desired whistle contour.

To guarantee that the generated positive samples have the same background magnitude as the input negative samples, we also include the identity loss,

$$L_{identity} = ||G_A(I_N, 0)||_1 + ||G_B(I_N) - (I_N, 0)||_1 \tag{5}$$

where 0 indicates an empty whistle contour input, i.e., we do not want the CycleGAN to generate any whistles. We denote $(I_N, 0)$ as the concatenated $I_N$ and empty whistle segmentation map. $G_A$ should produce residuals of zero when

there are no input whistle contours. We also use adversarial loss, $L_{D_A}, L_{D_B}, L_{G_A}, L_{G_B}$, and cycle consistence loss $(L_{cyc})$ from CycleGAN

$$L_{D_A} = (D_A(I_P) - 1)^2 + (D_A(I'_P))^2 \tag{6}$$

$$L_{D_B} = (D_B(I_N, \ I_W) - 1)^2 + (D_B(G_B(I_P)))^2 \tag{7}$$

$$L_{G_A} = (D_A(I'_P) \ - \ 1)^2 \tag{8}$$

$$L_{G_B} = (D_B(G_B(I_P)) - \ 1)^2 \tag{9}$$

$$L_{cyc} = ||G_B(I'_P) - (I_N, \ I_W)||_1 + \\ ||G_A(G_B(I_P)) - I_P||_1 \tag{10}$$

where $I_P$ refers to real positive samples. We simplify the notation of two $D_B$ networks in one $D_B$ function in the above equation. The full objective for generators is

$$L_G = L_{G_A} + L_{G_B} \ + \lambda_0 L_{cyc} + \lambda_1 L_{consistence} + \lambda_2 L_{identity} \tag{11}$$

where $\lambda_0$, $\lambda_1$, and $\lambda_2$ control the relative importance of the corresponding loss items. The full objective of the discriminator is

$$L_D = L_{D_A} + L_{D_B} \tag{12}$$

Ideally, $D_A$, $D_B$ will assign 1 to real samples and assign 0 to generated samples with this training objective. $G_A$, $G_B$ will try to fool the discriminators and generate realistic samples.

### C. Whistle extraction model

We use the whistle extraction model from [19] as our baseline. This model, which is similar to a selective edge detection model, produces a confidence map of the whistle signals. Although the generated samples are visually similar to real samples (Fig. 3), the distributions of the real and generated whistle contour may differ due to the imperfect training of GAN when data are limited. This discrepancy decreases the accuracy of our whistle extraction model when we use the generated samples for data augmentation. Therefore, we use ABN layers [31]; i.e., we use auxiliary BatchNorm (BN) layers for forwarding generated samples and normal BN layers for real samples. We share the same convolutional layers for real and generated samples. By denoting the input sample as $x$, the whistle signal label as $y$, and the whistle extraction model as $f$, the loss without ABN can be described as

$$L = ||y - f(x))||_2 \tag{13}$$

The loss with ABN is

$$L = \frac{1}{(1 + \lambda)} (||(y_{real} - f(x_{real}))||_2 + \\ \lambda ||y_{fake} - f_{abn}(x_{fake}))||_2) \tag{14}$$

where $x_{real}$, $y_{real}$ are the real samples and labels, respectively, and $x_{fake}$, $y_{fake}$ are the generated samples and labels, respectively. $\lambda$ is a factor to adjust the weights of real data and generated data in loss calculation. $f_{abn}(x)$ denotes the output of the whistle extraction model for input $x$ when the auxiliary BN layer is used in forwarding. We empirically find
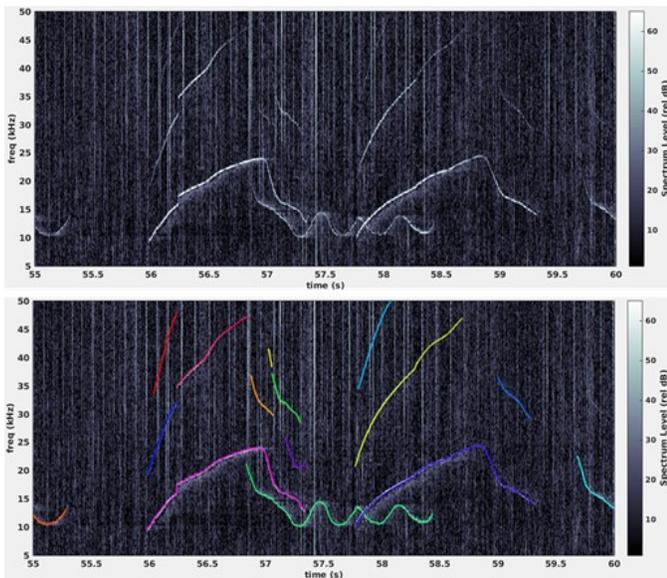
Fig. 4. Illustration of whistle extraction. (Top) spectrogram visualized by *Silbido* [11]; (Bottom) extracted whistles, where each whistle is highlighted with a different color.
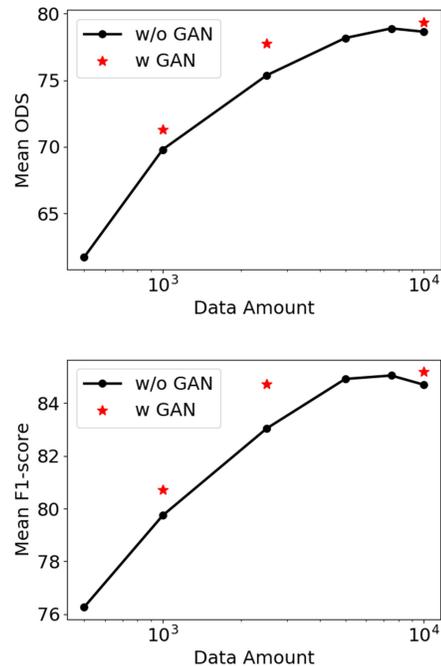


Fig. 5. Mean spectral peak detection F1-score (upper) or mean whistle extraction F1-score (lower) against the number of real positive samples in the training set. Optimal Dataset Scale (ODS) is an edge detection metric that assesses peak detection. "w/o GAN" and "w GAN" indicates the performance without and with GAN augmentation, respectively.

that ABN layers improve the whistle extraction performance when the distributions of the generated and real samples may be different.

The quality of GAN-synthesized samples is affected by the number of real samples available for training. The generator may synthesize poor-quality samples when the number of real examples used in GAN training is low. Fig. 3 provides examples of synthetic whistle contours when 2500 real positive samples are used for GAN training, including whistle contours that are of poor quality. Therefore, we designed two heuristic conditions for selecting high-quality generated samples. Denoting the value of an individual bin in the whistle contour patch as p, we select the sample for training the whistle extraction model when

$$\sum -plogp \ < T_e \tag{15}$$

and

$$\sum \delta(p - T_c) > T_p \tag{16}$$

where

$$\delta(x) = \begin{cases} 0 & x \le 0 \\ 1 & x > 0 \end{cases} \tag{17}$$

$T_e$ is a threshold for the sum of the pixel entropy, so the first condition removes generated whistles with diffuse medium-intensity signals (high entropy). The second condition chooses samples in which more than $T_p$ bins have intensity above $T_c$, allowing samples with short whistle fragments to be removed.

## IV. DATA AND IMPLEMENTATION

### A. Datasets

We used the whistle extraction data from the 2011 workshop on detection, classification, localization, and density estimation of marine mammals (DCLDE 2011, available on

the MobySound Archive [63]). These data contain recordings of calls made by five toothed whale species: long-beaked common dolphins (*Delphinus capensis*), short-beaked common dolphins (*Delphinus delphis*), bottlenose dolphins (*Tursiops truncatus*), melon-headed whales (*Peponocephala electra*), and spinner dolphins (*Stenella longirostris*). Whistle contours were annotated by trained analysts across the 5-50 kHz bandwidth as described in [11]. We use 30 recordings from the 5 species to train and 12 recordings from 4 species to test. Short-beaked common dolphins are removed from evaluation because some of the files had annotation errors. The training data consisted of approximately 127 min of recordings with 12,539 annotated whistles. The test data (∼43 min of acoustic data) contained 6,011 annotated whistles.

We computed log-magnitude spectrograms for the whistle extraction model and the GAN-based data synthesis. We employed series of discrete Fourier transforms in spectrogram computation. 8 ms Hamming-windowed frames (125 Hz bandwidth) were computed every 2 ms, and we empirically restricted the dynamic range of the $log_{10}$ magnitude spectrogram to the range [0, 6] (an intensity range of 0 to 120 dB rel.), i.e., we transformed the values <0 to 0, and those >6 to 6. We divided the spectrogram values by 6 which made them within [0, 1], and discarded the spectrogram values outside of the annotation frequency range of 5-50 kHz (361 frequency bins), which covers the frequency range of most delphinid whistles and their harmonics.

For network training, we partitioned the spectrogram into $64 \times 64$ patches, where each patch covered a time interval of

128 ms and frequency interval of 8 kHz. For the training data, we selected the positive patches with a sliding window with a 25 pixel step size across portions of spectrograms containing whistles, which led to 115,968 positive patches available for training. We randomly selected the same number of negative patches, which only contain noise, and combined them with positive patches as our training data (referred to as the full dataset). Most of our experiments used a subset of the full data (referred to as a reduced dataset). We describe the details of generating the reduced dataset in Section V-A.

### B. Networks and Algorithms

*1) Whistle extraction network:* We used the same network architecture as [19]. The model has 10 convolutional layers, including 1 input layer, 4 residual blocks (each block contains two convolution layers), and 1 output layer. The input layer and output layer use kernel size 5 and padding size of 2, and other layers use a kernel size of 3 and a padding size of 1. All hidden layers have 32 channels. The model input is a one-channel spectrogram and the output is a confidence map of whistle occurrence. The size of the output confidence map is the same as that of the input spectrogram.

We trained the whistle extraction model with an Adam optimizer (initial learning rate=$1 \times 10^{-3}$, betas = [0.9, 0.999], weight decay=$5 \times 10^{-4}$) for $1 \times 10^6$ and $3 \times 10^5$ iterations on the full dataset and reduced datasets, respectively. The learning rate was multiplied by 0.1 every $4 \times 10^5$ and $1 \times 10^5$ iterations for the full and reduced datasets, respectively. We set the batch size to 64, and we used 64 real samples and 64 generated samples in each iteration for data augmentation experiments. We used $\lambda$=1 in the loss function of Eq. 14 for our experiments with generated data, which make the generated samples have the same contribution of loss as real examples.

*2) WGAN:* We used the same WGAN architecture for the generation of whistle contours and negative samples. The generator network uses a fully-connected layer to output feature maps of size (512,4,4) from a 128-dimensional standard Gaussian distribution. Four groups of convolutional layers and pixel shuffle layers are used to gradually enlarge the feature map to $64 \times 64$. A Tanh layer is used to output the $64 \times 64$ patch. The discriminator network takes the generated samples and real samples as input, and outputs the Wasserstein distance estimation. It contains 4 convolutional layers with a stride of 2 and a fully connected layer. The networks are optimized by Adam optimizers (initial learning rate = $1 \times 10^{-4}$, betas = [0.5, 0.9], batch size = 64) for $3 \times 10^4$ and $5 \times 10^4$ iterations on the reduced and full datasets, respectively. In each WGAN training iteration, the discriminator is optimized for 5 steps while the generator is optimized for 1 step, where the network parameters are updated by applying the optimizer to one mini-batch of data in each step. For sample selection, we used $T_e$=70, $T_c$=0.5, $T_p$=64.

*3) CycleGAN:* The GAN model that we used to add whistles on synthetic noise employs the CycleGAN architecture of [30]. The generators follow the U-Net [64] architecture, which has 6 U-Net blocks with a basic width of 64. InstanceNorm layers are used in the U-Net blocks. The discriminator is a fully convolutional network with 3 convolutional layers. We trained the generators and discriminators with Adam optimizers (learning rate = $2 \times 10^{-4}$, betas = [0.5, 0.999], batch size = 64) for 25,120 iterations (160 epochs for 10,000 real positive samples) for the reduced dataset and 50 epochs for the full dataset. We set $\lambda_0$=10, $\lambda_1$=0.5, and $\lambda_2$=0.5 for Eq. 11. We apply a random $\gamma$ following a unified distribution between (0.5, 1.5) in Eq. 3.

*4) Graph Search:* We adapted the graph search [11] algorithm to the outputs of the whistle extraction network to predict individual whistles. This algorithm maintains sets of graphs, the nodes of which indicate the trace of predicted whistle contours. Multiple crossing whistles can be represented by a single graph. At each time step, local maximum points (peaks) on the confidence map are selected along the frequency dimension, and peaks with confidence greater than 0.5 are retained as candidate points. For each candidate point, the algorithm either initiates a new graph or extends terminating nodes of existing graphs. Extensions are made when the new node is along a reasonable trajectory predicted by a low-order polynomial fit of the graph path near a terminating node. Graphs that have not been extended within a specified time are considered closed. Closed graphs are removed from the current graph set. When a graph is of a shorter duration than a settable minimum whistle duration, it is discarded. Otherwise individual whistles are extracted from the graph on the basis of an analysis of graph vertices.

### C. Metrics

*1) Evaluation of confidence maps:* We first assessed the quality of the whistle-energy confidence maps predicted by the whistle-extraction model. To do this, we utilized the BSDS 500 benchmark tools and protocol [65] to calculate the highest dataset-scale F1-score across various thresholds (referred to as the "Optimal Dataset Scale," or ODS). We thinned each ground-truth whistle to a width of one pixel and compared them to predicted confidence maps that were binarized using 50 evenly distributed thresholds between 0 and 1. All default parameters within the benchmark tool were used in our evaluation.

*2) Evaluation of whistle extraction:* We used *Silbido* [11] to evaluate the quality of whistle extraction after the graph search was applied to the confidence map. This library calculates recall, the percentage of validated whistle contours that were detected; and precision, the percentage of detections that were correct. Then we calculated the precision, recall, and F1-score on testing files of each species and averaged them among all species. We determined the success or failure of whistle extraction results by examining the set of expected analyst annotations as described in [11]. We checked whether any of the detections overlapped with the analyst-annotated whistle contour in time. If so, we examined whether each overlapping detection matched the analyst's annotation. When the average deviation in frequency between the detected contour and annotation was < 350 Hz and the analyst detections had lengths ≥ 150 ms, with a signal-to-noise ratio ≥ 10 dB over at least 30% of the whistle, we classified the overlapping detections as matched detections. When an annotated whistle
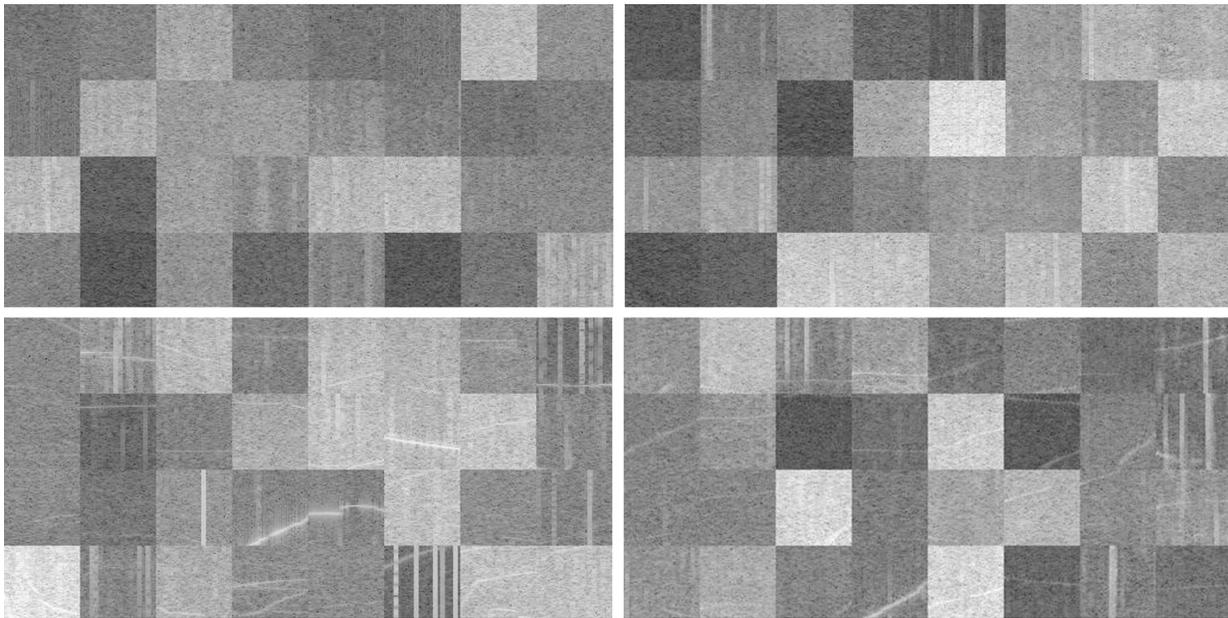
Fig. 6. Real background noise samples (upper left); Our GAN generated background noise samples (upper right); Real whistle samples (bottom left); Our GAN generated whistle samples (bottom right). Multiple 64 × 64 patches are concatenated in each category for better visualization of the data variance.

TABLE I
PERFORMANCE OF WHISTLE EXTRACTION

| N[a] | Mean ODS | | Mean F1-score | | Mean Precision | | Mean Recall | |
|---|---|---|---|---|---|---|---|---|
| | w/o GAN[b] | w GAN | w/o GAN | w GAN | w/o GAN | w GAN | w/o GAN | w GAN |
| 1000 | 69.80±2.41[c] | 71.33±2.58 | 79.74±2.94 | 80.73±1.89 | 85.01±3.54 | 76.86±3.99 | 84.82±3.44 | 78.32±3.58 |
| 2500 | 75.37±1.50 | 77.78±0.89 | 83.04±1.10 | 84.73±0.90 | 85.88±1.92 | 86.38±1.77 | 81.29±2.26 | 83.72±1.32 |
| 10000 | 78.64±0.67 | 79.38±0.38 | 84.70±1.11 | 85.21±0.80 | 87.55±1.52 | 87.13±1.86 | 82.67±1.52 | 83.85±1.13 |
| all | 80.85±0.23 | 81.23±0.10 | 87.42±0.44 | 87.88±0.14 | 89.27±0.20 | 89.60±0.31 | 86.04±0.67 | 86.63±0.36 |

[a] We denote the number of real positive samples for whistle extraction model and GAN training as N; "all" indicate that the full dataset is used.

[b] w GAN and w/o GAN indicate the performance of whistle extraction model with or without our GAN generated samples, respectively. The whistle extraction model is the same as [19].

[c] We conduct repeated experiments for each setting, and we report performance average ± standard deviation for each metric. Refer to Section IV-C for more details.

did not meet the above criteria (too short or low intensity), we discarded any matching detections, and they did not contribute to the metrics. We classified unmatched detections as false positives.

## V. EXPERIMENTS AND RESULTS

### A. *Varied number of annotated samples*

We first studied the effect of varying the amount of training data for our whistle extraction network. Because annotation is expensive, a key motivation for data augmentation is to reduce the number of annotations required. Training effective deep-learning models requires a considerable amount of high-quality annotated data [66]. For the whistle extraction task in this paper, it remains unclear how the whistle models perform when the amount of annotated data varies. To address this issue, we conducted 6 experiments that selected n positive patches and n negative patches, where n = 500, 1000, 2500, 5000, 7500, or 10000. Random selection of patches was structured to ensure that smaller datasets were subsets of larger ones. We repeated this process five times to obtain 5 datasets for each n. For each dataset, we trained whistle extraction models 5 times, and report average performance.

The experimental results are shown in Fig. 5. The black curves show the performance of the confidence map (ODS) and whistle extraction (F1-score) (upper and lower plots, respectively) with respect to the quantity of training data. While the ODS quantifies the performance of the whistle extraction model in detecting the presence and shape of the whistles, the results suggest that, with more training data, the average ODS increases. The increase in whistle extraction F1-score follows the same trend as ODS. Our results show that increasing the amount of annotated data substantially improves the performance of whistle detection. At the same time, as the amount of data increases, the rate of performance improvement decreases, which means that exponentially more data may be needed to increase performance by 1 unit when the initial dataset is larger.

### B. *Data augmentation*

We also studied the effect of varying dataset size on GAN training and data augmentation. In this set of experiments, we applied the proposed augmentation method to augment n = 1000, 2500, and 10000 positive samples and negative samples. In each experiment, we generated 10 × n samples with our
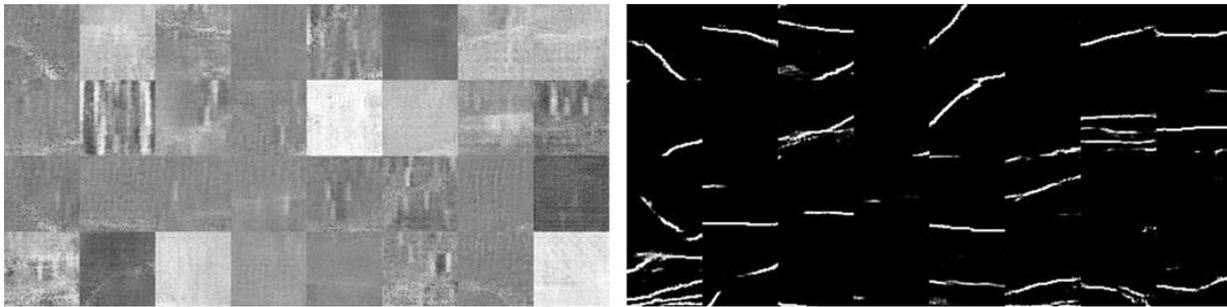
Fig. 7. Positive samples (left) and corresponding whistle contour (right) generated by vanilla GAN. Multiple $64 \times 64$ patches are concatenated.

WAS-GAN. All GAN networks were randomly initialized and trained once per dataset. For each augmented dataset, we trained the whistle extraction model with ABN for 5 times.

Fig. 6 shows examples of samples generated by our WAS-GAN (n = 2500). By visually comparing the real samples and generated samples, we see that the noise patterns and whistle signal patterns are well simulated by our GAN networks, e.g., the clicks (wide vertical band of high energy across the frequency domain) are simulated well, as are the width and strength of whistle signals.

Table I reports the experiment's ability to correctly predict time-frequency peaks associated with whistles (mean ODS) and to correctly extract whistles from these predictions (mean F1-score). Consistent performance improvements were obtained for both measures. Our methods obtained gains of 1.53, 2.41, and 0.74 in mean ODS, and 0.99, 1.69, and 0.51 in mean F1-score for the three augmentation experiments when n=1000, 2500, and 10000 training patches, respectively. We also obtained improvements of 0.38 and 0.46 in the mean ODS and mean F1-score, respectively, when we used WAS-GAN on the full dataset. In comparison to experiments using n=10000, we utilized over 100,000 additional annotated samples in our full dataset experiment. These samples were manually labeled as opposed to our GAN augmented samples, and this led to an increase of 2.72 in the whistle extraction F1-score. Without our GAN-generated samples, in order to achieve a 0.46 increase in the F1-score by adding more human-annotated samples to our current dataset, we would have to annotate tens of thousands more samples. The training stability was notably improved (with a reduction in the variance of the F1 metrics) with the addition of the generated data. These improvements highlight the effectiveness of our proposed stage-wise, GAN-based data augmentation method: the use of augmented data improves spectral peak detection results, which in turn also improves whistle contour extraction results.

## C. Ablation study

We conducted a set of ablation experiments to examine the contributions of different components of the proposed method. We chose datasets with n = 2500 samples for these experiments. The quantitative results are shown in Table III.

*1) Residual learning:* In this ablation experiment, we trained the CycleGAN in stage 3 to directly generate positive samples rather than adding the residual to the negative samples

(Eq. 3). While we can change the whistle signal magnitude by altering the weight in Eq. 3 when the generator outputs residual, the whistle signal's magnitude is determined by the generator model in this setting. In contrast to the proposed WAS-GAN, we observed a decrease of 1.43 in mean ODS and a decrease of 1.44 in mean whistle extraction F1-score when we removed residual learning. This performance drop might be caused by the fact that the GAN needed to output background noise, which might increase the difficulty and instability of learning. Moreover, the variance of generated data decreases when the magnitude of whistle signals cannot be adjusted by the multiplier of the residual.

*2) Patch selection:* This ablation experiment removed the quality assurance filter (Eq. 15 and Eq. 16) for whistles generated by the GAN. As a result, generated whistles similar to those surrounded by the red bounding boxes (Fig. 3) were included in the training data. The mean ODS dropped by 1.21 and the mean F1-score decreased by 1.44 after this change. This indicates that low-quality samples may reduce the performance of the whistle extraction network training, and our simple heuristic selection method effectively selects samples for the whistle extraction task.

TABLE II
ABLATION STUDY

| Experiments | Mean ODS | Mean F1-score |
|---|---|---|
| 2500+GAN[a] | 77.78 | 84.73 |
| - residual[b] | -1.43 | -1.44 |
| - select | -1.21 | -1.44 |
| - ABN | -0.68 | -0.98 |
| - ABN, - select | -0.86 | -1.97 |
| - residual, - select, - ABN | -2.01 | -5.21 |
| vanilla GAN[c] | -0.57 | -1.04 |
| Random Addition[d] | -0.36 | -0.65 |
| Random Addtion + Gaussian Blur[e] | -0.37 | -0.67 |

[a] GAN augmentation from 2500 real positive samples and 2500 negative samples. We report the whistle extraction performance with our proposed GAN method in this row and the change of performance compared to this row in the following rows.

[b] -XXX means that component XXX is removed. The components include: (i) residual: residual learning; (ii) select: selection of synthetic whistles with entropy and duration criteria; (iii) ABN: auxiliary batch normalization.

[c] We replace stage-wise GANs with a single WGAN-gp [48] for sample synthesis.

[d] We remove the third GAN model (CycleGAN) and directly add the output of the first two GANs with random weights.

[e] We apply random Gaussian blurring to the generated whistle contour before it is added to background noise.
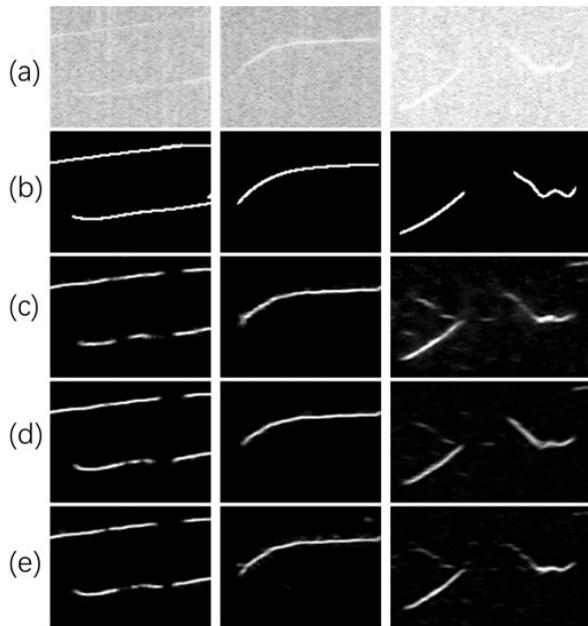
Fig. 8. Outputs of whistle extraction models. Models with the best whistle extraction F1-score among all parallel experiments in each training setting are visualized. (a) Spectrograms that are used as model input. (b) Ground truth. (c) Output of model trained with 2500 real positive patches and negative patches. (d) Output of model trained with 2500 positive patches and negative patches and GAN synthesized data. (e) Same as d, but the model does not have auxiliary batch normalization (ABN).

*3) ABN:* Because ABN stores statistics of real samples and generated samples separately, it may better stabilize the training when the generated samples and real samples have different distributions [31]. We evaluated the functionality of applying ABN with and without patch selection to our whistle extraction task; patch selection affects the generated sample distribution. After removal of ABN, the whistle extraction F1-score dropped by 0.98 with patch selection and by 1.97 without patch selection. This suggests that our patch selection method contributes to generating samples that are closer to the actual distribution of whistles. The performance change is consistent with our hypothesis that generated samples and real samples have a different distribution when few data are included in GAN training.

We also observed decreases in ODS of 0.68 with patch selection and 0.86 without patch selection after removal of ABN, which is a less decrease compared to whistle extraction F1-score. While ODS demonstrates the whistle extraction model's performance at the spectrogram bin level, this metric does not always linearly correlate to the whistle extraction performance, because it ignores the signal continuity among bins. We observed that removing ABN frequently resulted in poorer continuity of predicted patches (e.g., Fig. 8d and 8e, first and third examples) and a greater number of false positives (e.g., Fig. 8d and 8e, second example). The whistle extraction F1-score also indicates the model's ability to recognize whistle signals under varying noise conditions or suppress false positives in the high-energy region of spectrogram according to the context information (signals in the neighborhood). The generated whistle contour and signals may be less continuous

than the real samples, which will train the whistle extraction model to ignore context information and make discontinuous predictions when ABN is removed. The comparison among Fig. 8 (c), (d), (e) rightmost column also shows that use of our generated data reduces false positives.

*4) Stage-wise GAN:* Instead of decomposing the sample generation into multiple stages, we used a single WGAN-gp with two output channels to generate whistle data, the spectrogram samples and their labels, similar to [53]. To deal with the increased learning difficulty of one WGAN, we increased the WGAN-gp capacity of the generator by using twice the number of hidden layers for each convolutional layer output as that in Section IV-B2. Examples of samples generated by this model are shown in Fig. 7. We saw clear artifacts and unnatural, sudden changes in the magnitude in adjacent bins on the spectrogram. The visual quality of generated samples was substantially worse than those generated by our stage-wise GAN in Fig. 6. We also observed a decrease of 1.04 in the whistle extraction F1-score compared to our proposed framework. Data augmentation with the low-quality samples still permitted the performance of the model to surpass that without augmentation for the time-frequency detection task. The negative effect of using corrupted data might be mitigated by the ABN layer.

*5) The third GAN:* In this ablation study, we remove the third GAN and instead generate positive sample $I'_P$ by simply adding the generated whistle contour $I_W$ to the generated background noise $I_N$. Following the work of Li et al. [19], we apply Gaussian blur $G$ with random deviation parameter $\sigma$ to the whistle contour, and we add the blurred contour to the background noise:

$$I'_P = I_N + \lambda CLIP(I_W + G(Y, \sigma)) \qquad (18)$$

where the clipping function $CLIP(x)$ is

$$CLIP(x) = \begin{cases} 0, & x \in (-\infty, 0) \\ x, & x \in [0, 1] \\ 1, & x \in (1, +\infty) \end{cases} \qquad (19)$$

We also try a simple version which does not contain Gaussian blur:

$$I'_P = I_N + \lambda I_W \qquad (20)$$

where $\lambda$ is a random weighting parameter. We use the same parameter setting as Li et al. [19], where $\lambda$ and $\sigma$ are uniform random numbers within the ranges of $[0.03, 0.23]$ and $[0.3, 1.3]$, respectively. As shown in Table III, both methods in Equation 19 and Equation 20 lead to inferior performance compared to the proposed stage-wise GAN method ("2500+GAN") that uses the third GAN. Considering that we use the same set of background noise and whistle contour shapes, this ablation study indicates that our proposed stage-wise GAN method generates more realistic whistle signals with a appearance which contributes to the improved training of the whistle model.

### D. Comparison with other whistle extraction methods

In addition to our previous work on network-based whistle extraction [19], we have selected two representative and

TABLE III
COMPARISON OF WHISTLE EXTRACTION METHODS

| Method | F1-score | Precision | Recall |
|---|---|---|---|
| Roch et al., 2011 [11] | 75.95 | 81.125 | 72.275 |
| Gruden et al., 2020 [22] | 83.40 | 76.55 | 92.45 |
| Gruden et al., 2020 ($\geq$150ms) | 74.38 | 95.85 | 60.875 |
| Li et al., 2020 [19] | 87.42 | 89.27 | 86.04 |
| Li et al., 2020 + our GAN | 87.88 | 89.60 | 86.63 |

competitive whistle extraction methods for comparison. Both methods identify whistle candidate points by determining if the Signal-to-Noise Ratio (SNR) values are above a threshold on the denoised spectrogram. The Graph-Search method developed by Roch et al. [11] employs graphs consisting of candidate points, which are extended with new points based on how well these new candidate points align with the existing graph through polynomial fitting results. As a point of comparison, Gruden et al. [22] uses a probabilistic approach based on the sequential Monte-Carlo probability hypothesis density (SMCPHD). In addition to the result of all SMCPHD predictions, we also present the results of predictions that are longer than 150ms, as both Graph-Search and our method apply this length criterion for detection.

Our approach outperforms SMCPHD and Graph-Search in the whistle extraction F1-score by 4.48 and 11.93, respectively. Additionally, our GAN-generated samples improve the method in [19] by 0.46 in F1-score, 0.33 in precision, and 0.59 in recall. SMCPHD demonstrates the highest recall but the lowest precision in this comparison, which indicates its aggressive strategy of making more whistle predictions. By removing whistle detections by SMCPHD that are shorter than 150ms, the precision of SMCPHD is improved by 19.3, while the recall is decreased by 31.57. This study suggests that SM-CPHD prefers shorter segments of whistles in its predictions. Our GAN-generated samples help the learning-based model achieve a competitive performance advantage on this whistle extraction task, however, it should be noted that optimizing the other algorithms for this specific dataset may diminish these advantages.

## VI. CONCLUSION AND DISCUSSION

We present a framework of stage-wise generative adversarial networks to generate training samples for whistle extraction. The data generation process consists of three stages: (i) generate time x frequency spectrogram patches containing background noise (ii) generate whistle contours and automatically discard poor quality contours (iii) fuse whistle signals with the background noise. Each stage is completed by one trained generative adversarial network. Compared to using a single vanilla GAN generating whistle extraction data and labels, our stage-wise GANs can generate samples with fewer artifacts which results in increased whistle extraction performance. We examined our data generation method by a series of experiments employing differing quantities of real and generated data, and note that using the generated data lead to consistent performance gains.

The stage-wise design may mainly contribute to the success of our data generation method. It separates the modeling of different components and the relationship between components, which eases the learning of the GANs in each stage as well as provides a straightforward way to explore different combinations of components. In our case, we generated the background noise separately and we were able to add different whistle signals to the same background. If we directly apply this idea to semantic segmentation data generation of natural images, we may first generate the appearance of background scene, then generate objects on it according to a desired segmentation map. If we extend this idea, we may generate the appearance of different objects separately and then add them to the background. In this way, we may fully explore combinations of varying objects and background appearances in the same segmentation layout. In our whistle extraction experiments, we did not use this extended idea, because the appearance of our foreground object (the whistles) is relatively simple, i.e., the variance of appearance is mainly rooted in the whistle contour shape and whistle magnitude. Therefore, we directly add whistle signals to the background using the third GAN in our framework. Our framework can be readily extended to extract calls of other whale species (e.g, baleen whales) and to other similar tasks (e.g., semantic image segmentation).

Though it may not affect the main contributions of this work, our data generation method can be improved in three aspects in the future. Firstly, we may use improved generative neural network architecture and training strategies. For example, we may use a generator architecture based on a style-transfer network which improves the generated sample quality [67]. The discriminator augmentation mechanism proposed in [68] may help stabilize training in limited data regimes. We may also explore generating larger patches of high quality with the method in [69]. Secondly, we may use real data in the data generation process to enrich the data variance. The real background and annotated whistle contours can be used as the input data of our GAN in the third stage, and we can generate whistle signsals of novel shapes on real background or generate whistle signals of annotated contour shapes on GAN-generated background. Thirdly, we may improve the sample selection method. In this paper, we use a simple yet effective pixel-wise entropy method to select whistle contour of good quality. Metric measuring texture or semantic information like [70] may better measure the quality of our generated samples and improve the sample selection process.

## REFERENCES

[1] L. Rabiner and B.-H. Juang, *Fundamentals of Speech Recognition*. Prentice-Hall, Inc., 1993.

[2] J. Ren, X. Jiang, J. Yuan, and N. Magnenat-Thalmann, "Sound-event classification using robust texture features for robot hearing," *IEEE Transactions on Multimedia*, vol. 19, no. 3, pp. 447–458, 2016.

[3] S. Kahl, C. M. Wood, M. Eibl, and H. Klinck, "BirdNET: A deep learning solution for avian diversity monitoring," *Ecological Informatics*, vol. 61, p. 101236, 2021.

[4] C.-H. Lee, J.-L. Shih, K.-M. Yu, and H.-S. Lin, "Automatic music genre classification based on modulation spectral analysis of spectral and cepstral features," *IEEE Transactions on Multimedia*, vol. 11, no. 4, pp. 670–682, 2009.

[5] A. Rizzi, M. Antonelli, and M. Luzi, "Instrument learning and sparse nmd for automatic polyphonic music transcription," *IEEE Transactions on Multimedia*, vol. 19, no. 7, pp. 1405–1415, 2017.

[6] S. Zhang, S. Zhang, T. Huang, and W. Gao, "Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching," *IEEE Transactions on Multimedia*, vol. 20, no. 6, pp. 1576–1590, 2017.

[7] W. A. Yost, "Fundamentals of hearing: An introduction," 2001.

[8] K. Vijayan, X. Gao, and H. Li, "Analysis of speech and singing signals for temporal alignment," in *2018 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2018, pp. 1893–1898.

[9] B. Met-Montot, S. Cabon, G. Carrault, and F. Porée, "Spectrogram-based fundamental frequency tracking of spontaneous cries in preterm newborns," in *2020 28th European Signal Processing Conference (EU-SIPCO)*. IEEE, 2021, pp. 1185–1189.

[10] W. T. Lu, L. Su *et al.*, "Vocal melody extraction with semantic segmentation and audio-symbolic domain transfer learning." in *International Society for Music Information Retrieval (ISMIR)*, 2018, pp. 521–528.

[11] M. A. Roch, T. Scott Brandes, B. Patel, Y. Barkley, S. Baumann-Pickering, and M. S. Soldevilla, "Automated extraction of odontocete whistle contours," *The Journal of the Acoustical Society of America*, vol. 130, no. 4, pp. 2212–2223, 2011.

[12] D. Gillespie, M. Caillat, J. Gordon, and P. White, "Automatic detection and classification of odontocete whistles," *The Journal of the Acoustical Society of America*, vol. 134, no. 3, pp. 2427–2437, 2013.

[13] J.-j. Jiang, L.-r. Bu, F.-j. Duan, X.-q. Wang, W. Liu, Z.-b. Sun, and C.-y. Li, "Whistle detection and classification for whales based on convolutional neural networks," *Applied Acoustics*, vol. 150, pp. 169–178, 2019.

[14] V. M. Janik, S. L. King, L. S. Sayigh, and R. S. Wells, "Identifying signature whistles from recordings of groups of unrestrained bottlenose dolphins (*Tursiops truncatus*)," *Marine Mammal Science*, vol. 29, no. 1, pp. 109–122, 2013.

[15] A. G. Taruski, "The whistle repertoire of the north atlantic pilot whale (*Globicephala melaena*) and its relationship to behavior and environment," in *Behavior of marine animals*. Springer, 1979, pp. 345–368.

[16] B. L. Sjare and T. G. Smith, "The relationship between behavioral activity and underwater vocalizations of the white whale, *Delphinapterus leucas*," *Canadian Journal of Zoology*, vol. 64, no. 12, pp. 2824–2831, 1986.

[17] F. Thomsen, D. Franck, and J. K. Ford, "On the communicative significance of whistles in wild killer whales (*Orcinus orca*)," *Naturwissenschaften*, vol. 89, no. 9, pp. 404–407, 2002.

[18] A. Jaramillo-Legorreta, G. Cardenas-Hinojosa, E. Nieto-Garcia, L. Rojas-Bracho, J. Ver Hoef, J. Moore, N. Tregenza, J. Barlow, T. Gerrodette, L. Thomas *et al.*, "Passive acoustic monitoring of the decline of mexico's critically endangered vaquita," *Conservation Biology*, vol. 31, no. 1, pp. 183–191, 2017.

[19] P. Li, X. Liu, K. Palmer, E. Fleishman, D. Gillespie, E.-M. Nosal, Y. Shiu, H. Klinck, D. Cholewiak, T. Helble *et al.*, "Learning deep models from synthetic data for extracting dolphin whistle contours," in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–10.

[20] A. Mallawaarachchi, S. Ong, M. Chitre, and E. Taylor, "Spectrogram denoising and automated extraction of the fundamental frequency variation of dolphin whistles," *The Journal of the Acoustical Society of America*, vol. 124, no. 2, pp. 1159–1170, 2008.

[21] H. Chen, J. Yan, N. U. R. Junejo, J. Qi, and H. Sun, "Sparse representation based on tunable q-factor wavelet transform for whale click and whistle extraction," *Shock and Vibration*, vol. 2018, 2018.

[22] P. Gruden and P. R. White, "Automated extraction of dolphin whistles—a sequential Monte Carlo probability hypothesis density approach," *The Journal of the Acoustical Society of America*, vol. 148, no. 5, pp. 3014–3026, 2020.

[23] X. Wang, J. Jiang, F. Duan, C. Liang, C. Li, Z. Sun, R. Lu, F. Li, J. Xu, and X. Fu, "A method for enhancement and automated extraction and tracing of *odontoceti* whistle signals base on time-frequency spectrogram," *Applied Acoustics*, vol. 176, p. 107698, 2021.

[24] J. Jeong, S. Lee, J. Kim, and N. Kwak, "Consistency-based semi-supervised learning for object detection," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[25] P. Li and X. Liu, "Learning knowledge-rich sequential model for planar homography estimation in aerial video," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 10584–10591.

[26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[27] T. Dao, A. Gu, A. Ratner, V. Smith, C. De Sa, and C. Ré, "A kernel theory of modern data augmentation," in *International Conference on Machine Learning*. PMLR, 2019, pp. 1528–1537.

[28] H. Brumm and S. A. Zollinger, "The evolution of the Lombard effect: 100 years of psychoacoustic research," *Behaviour*, vol. 148, no. 11-13, pp. 1173–1198, 2011.

[29] W. W. Au, D. A. Carder, R. H. Penner, and B. L. Scronce, "Demonstration of adaptation in beluga whale echolocation signals," *The Journal of the Acoustical Society of America*, vol. 77, no. 2, pp. 726–730, 1985.

[30] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.

[31] C. Xie, M. Tan, B. Gong, J. Wang, A. L. Yuille, and Q. V. Le, "Adversarial examples improve image recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 819–828.

[32] A. Wali, Z. Alamgir, S. Karim, A. Fawaz, M. B. Ali, M. Adan, and M. Mujtaba, "Generative adversarial networks for speech processing: A review," *Computer Speech & Language*, vol. 72, p. 101308, 2022.

[33] F. Dadouchi, C. Gervaise, C. Ioana, J. Huillery, and J. I. Mars, "Automated segmentation of linear time-frequency representations of marine-mammal sounds," *The Journal of the Acoustical Society of America*, vol. 134, no. 3, pp. 2546–2555, 2013.

[34] X. C. Halkias and D. P. Ellis, "Call detection and extraction using bayesian inference," *Applied Acoustics*, vol. 67, no. 11-12, pp. 1164–1174, 2006.

[35] P. White and M. Hadley, "Introduction to particle filters for tracking applications in the passive acoustic monitoring of cetaceans," *Canadian Acoustics*, vol. 36, no. 1, pp. 146–152, 2008.

[36] P. Gruden and P. R. White, "Automated tracking of dolphin whistles using gaussian mixture probability hypothesis density filters," *The Journal of the Acoustical Society of America*, vol. 140, no. 3, pp. 1981–1991, 2016.

[37] D. K. Mellinger, S. W. Martin, R. P. Morrissey, L. Thomas, and J. J. Yosco, "A method for detecting whistles, moans, and other frequency contour sounds," *The Journal of the Acoustical Society of America*, vol. 129, no. 6, pp. 4055–4061, 2011.

[38] A. Kershenbaum and M. A. Roch, "An image processing based paradigm for the extraction of tonal sounds in cetacean communications," *The Journal of the Acoustical Society of America*, vol. 134, no. 6, pp. 4435–4445, 2013.

[39] O. Serra, F. Martins, and L. R. Padovese, "Active contour-based detection of estuarine dolphin whistles in spectrogram images," *Ecological Informatics*, vol. 55, p. 101036, 2020.

[40] K. Han and D. Wang, "Neural network based pitch tracking in very noisy speech," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 12, pp. 2158–2168, 2014.

[41] R. M. Bittner, B. McFee, J. Salamon, P. Li, and J. P. Bello, "Deep salience representations for $F_0$ estimation in polyphonic music." in *International Society for Music Information Retrieval (ISMIR)*, 2017, pp. 63–70.

[42] S. Liu, M. Liu, M. Wang, T. Ma, and X. Qing, "Classification of cetacean whistles based on convolutional neural network," in *2018 10th International Conference on Wireless Communications and Signal Processing (WCSP)*. IEEE, 2018, pp. 1–5.

[43] M. Ngo, S. Karaoglu, and T. Gevers, "Self-supervised face image manipulation by conditioning GAN on face decomposition," *IEEE Transactions on Multimedia*, 2021.

[44] L. Galteri, L. Seidenari, M. Bertini, and A. Del Bimbo, "Deep universal generative adversarial compression artifact removal," *IEEE Transactions on Multimedia*, vol. 21, no. 8, pp. 2131–2145, 2019.

[45] J. Lu, W. Zhang, and H. Yin, "Generate and purify: Efficient person data generation for re-identification," *IEEE Transactions on Multimedia*, 2021.

[46] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in Neural Information Processing Systems*, vol. 27, 2014.

[47] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[48] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of Wasserstein GANs," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[49] H. Emami, M. M. Aliabadi, M. Dong, and R. B. Chinnam, "SPA-GAN: Spatial attention GAN for image-to-image translation," *IEEE Transactions on Multimedia*, vol. 23, pp. 391–401, 2020.

[50] L. Chen, L. Wu, Z. Hu, and M. Wang, "Quality-aware unpaired image-to-image translation," *IEEE Transactions on Multimedia*, vol. 21, no. 10, pp. 2664–2674, 2019.

[51] M. Frid-Adar, E. Klang, M. Amitai, J. Goldberger, and H. Greenspan, "Synthetic data augmentation using GAN for improved liver lesion classification," in *2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018)*. IEEE, 2018, pp. 289–293.

[52] G. Mariani, F. Scheidegger, R. Istrate, C. Bekas, and C. Malossi, "BAGAN: Data augmentation with balancing GAN," *arXiv preprint arXiv:1803.09655*, 2018.

[53] C. Bowles, L. Chen, R. Guerrero, P. Bentley, R. Gunn, A. Hammers, D. A. Dickie, M. V. Hernández, J. Wardlaw, and D. Rueckert, "GAN augmentation: Augmenting training data using generative adversarial networks," *arXiv preprint arXiv:1810.10863*, 2018.

[54] S.-W. Huang, C.-T. Lin, S.-P. Chen, Y.-Y. Wu, P.-H. Hsu, and S.-H. Lai, "AugGAN: Cross domain adaptation with GAN-based data augmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 718–731.

[55] J. Zeng, Q. Chen, Y. Liu, M. Wang, and Y. Yao, "StrokeGAN: Reducing mode collapse in chinese font generation via stroke encoding," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 4, 2021, pp. 3270–3277.

[56] T. Osakabe, M. Tanaka, Y. Kinoshita, and H. Kiya, "CycleGAN without checkerboard artifacts for counter-forensics of fake-image detection," in *International Workshop on Advanced Imaging Technology (IWAIT) 2021*, vol. 11766. International Society for Optics and Photonics, 2021, p. 1176609.

[57] S. Mun, S. Park, D. K. Han, and H. Ko, "Generative adversarial network based acoustic scene training set augmentation and selection using SVM hyper-plane," *Proceedings of the Detection and Classification of Acoustic Scenes and Events (DCASE)*, pp. 93–97, 2017.

[58] S. Agarwal and H. Farid, "Detecting deep-fake videos from aural and oral dynamics," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 981–989.

[59] A. Antoniou, A. Storkey, and H. Edwards, "Data augmentation generative adversarial networks," *arXiv preprint arXiv:1711.04340*, 2017.

[60] J. Mu, W. Qiu, G. D. Hager, and A. L. Yuille, "Learning from synthetic animals," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 386–12 395.

[61] A. Waheed, M. Goyal, D. Gupta, A. Khanna, F. Al-Turjman, and P. R. Pinheiro, "CovidGAN: data augmentation using auxiliary classifier GAN for improved COVID-19 detection," *IEEE Access*, vol. 8, pp. 91 916–91 923, 2020.

[62] S. Pandey, P. R. Singh, and J. Tian, "An image augmentation approach using two-stage generative adversarial network for nuclei image segmentation," *Biomedical Signal Processing and Control*, vol. 57, p. 101782, 2020.

[63] D. K. Mellinger and C. W. Clark, "Mobysound: A reference archive for studying automatic recognition of marine mammal sounds," *Applied Acoustics*, vol. 67, no. 11-12, pp. 1226–1242, 2006, (Webpage last viewed on May 23, 2022). [Online]. Available: http://www.mobysound.org/workshops_p2.html

[64] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[65] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 5, pp. 898–916, 2010.

[66] A. Althnian, D. AlSaeed, H. Al-Baity, A. Samha, A. B. Dris, N. Alzakari, A. Abou Elwafa, and H. Kurdi, "Impact of dataset size on classification performance: an empirical evaluation in the medical domain," *Applied Sciences*, vol. 11, no. 2, p. 796, 2021.

[67] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, 2019, pp. 4401–4410.

[68] T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, and T. Aila, "Training generative adversarial networks with limited data," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 104–12 114, 2020.

[69] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of GANs for improved quality, stability, and variation," in *International Conference on Learning Representations*, 2018.

[70] K. G. Larkin, "Reflections on Shannon Information: In search of a natural information-entropy for images," *arXiv preprint arXiv:1609.01117*, 2016.