

# Data Augmentation for Environmental Sound Classification Using Diffusion Probabilistic Model with Top-k Selection Discriminator

Yunhao Chen<sup>1</sup>[0000-0002-8134-2314], Yunjie Zhu<sup>2</sup>, Zihui Yan<sup>1</sup>, and Jianlu Shen<sup>1</sup>  
Zhen Ren<sup>1</sup> Yifan Huang<sup>1</sup>

<sup>1</sup> Jiangnan University, Wuxi, 214000, China  
1191200221@stu.jiangnan.edu.cn

<sup>2</sup> University of Leeds, Leeds, LS2 9JT, United Kingdom

**Abstract.** Despite consistent advancement in powerful deep learning techniques in recent years, large amounts of training data are still necessary for the models to avoid overfitting. Synthetic datasets using generative adversarial networks (GAN) have recently been generated to overcome this problem. Nevertheless, despite advancements, GAN-based methods are usually hard to train or fail to generate high-quality data samples. In this paper, we propose an environmental sound classification augmentation technique based on the diffusion probabilistic model with DPM-Solver++ for fast sampling. In addition, to ensure the quality of the generated spectrograms, we propose a top-k selection technique to filter out the low-quality synthetic data samples. According to the experiment results, the synthetic data samples have similar features to the original dataset and can significantly increase the classification accuracy of different state-of-the-art models compared with traditional data augmentation techniques. The public code is available on <https://github.com/JNAIC/DPMs-for-Audio-Data-Augmentation>.

**Keywords:** diffusion probabilistic models, data augmentation, environmental sound classification

## 1 Introduction

In recent years, deep learning models have been advancing in the sound classification task like the convolutional neural networks (CNN)[2], [6], [7], transformer-based network [3], [4], CNN-RNN network[5]. However, these methods are hungry for considerable amounts of data for efficient performance due to their large amounts of parameters. Consequently, the main challenge the deep learning methods are confronted with is the limited samples of the training dataset. The most challenging aspect of developing deep learning supervised models is data annotation. It is a labour-intensive process that is expensive and time-consuming, especially when the deep learning methods require a lot of samples. Consequently, data augmentation is proposed to overcome the limited data samples problem.

Traditional data augmentation for audio classification includes flip, rotation, scale, crop, translation, noise, pitch shifting, masking, etc.[8,2] However, these methods are based on simple or linear transformations. Compared with the complexity of deep learning-based methods, these augmentations can not efficiently improve the performance of CNN-based or transformer-based classifier methods. As a result, to enhance the accuracy of deep learning-based methods, it is necessary to either apply data augmentation techniques that involve transformations of similar complexity to the deep learning models, or employ a model that can represent the probability distribution of the real data samples. Accordingly, researchers take advantage of the generative models, especially the generative adversarial network (GAN) [8,2,5,9,10] to synthesize new data samples with qualities similar to real data samples.

However, despite the considerable amounts of GANs’ applications in data augmentation, they are subject to unstable training processes, model collapse issues and failure to represent a broad enough data distribution. [11,12,13,14] Consequently, the GANs are challenging to be scaled rapidly to a new application.

On the other hand, diffusion probabilistic models [15,16] are becoming much more popular than GANs. Moreover, the diffusion-based generative models have been proven that they can beat the GANs on many tasks such as image synthesis [17], medical image synthesis where the data samples are limited[18], topology optimization[19] and so on. Therefore, we adopt diffusion probabilistic models for high-quality data augmentation. However, the popular diffusion models’ sampling procedures are mainly based on denoising diffusion implicit models (DDIM) [20], which require 50 to 100 steps to generate high-quality data samples. This is extremely time-consuming when generating thousands of data samples. On the other hand, DPM-Solver++ [21] only requires 10 to 20 steps to generate similar results compared with the DDIM method. Consequently, we employ the DPM-Solver++ for our sampling strategy.

Diffusion probabilistic models can generate diverse data samples from complex distributions by reversing a Markov chain of Gaussian diffusion processes. However, the quality of the data samples generated by these models is not always satisfactory, as they may contain artefacts, blurriness or inconsistency with the target distribution. Consequently, we propose a Top-k Selection method to filter out inappropriate data samples based on a pretrained model to address this problem. Our method can improve the quality of data generation without modifying the diffusion probabilistic models.

This paper aims to improve the environmental sound classification process using a conditional diffusion probabilistic network framework for high-quality data augmentation with a discriminator and DPM-Solver++. To summarize, the main contributions of this paper are as follows:

- 1) We present the first study on applying conditional diffusion probabilistic network frameworks to generate high-quality data samples for environmental sound classification tasks based on a popular sound dataset, UrbanSound8K [22].

2) We propose a post-processing approach called top-k selection based on a pre-trained discriminator. This approach automatically eliminates samples with low quality and insufficient representation after training.

3) We evaluate seven state-of-the-art deep learning (DL) models for environmental sound classification on the dataset with real and synthetic data samples generated by data augmentation. We train the models from scratch, without transfer learning, and show significant accuracy improvement with synthetic data.

## 2 Method

This section explains the methods for diffusion probabilistic models, DPM-Solver++ and data augmentation.

### 2.1 Diffusion Probabilistic Models

Diffusion probabilistic models (DPMs) [23][24] are a class of generative models that convert Gaussian noise into samples from a learned data distribution via an iterative denoising process. Non-equilibrium thermodynamics serves as the basis for diffusion models. To gradually introduce random noise to the data, [23][24] establish a Markov chain of diffusion steps. Then they figure out how to reverse the diffusion process to create the desired data samples from the noise by virtual deep learning methods, the details of which will be discussed next.

**Forward Diffusion Process** A forward diffusion process is defined as a process that adds Gaussian noise to a data sample  $\mathbf{x}_0$  sampled from a real data distribution  $q(\mathbf{x})$  over  $T$  steps, resulting in a sequence of noisy samples  $\mathbf{x}_1, \dots, \mathbf{x}_T$ . The amount of noise added at each step is determined by a variance schedule  $(\beta_t \in (0, 1)_{t=1}^T)$ .

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (1)$$

The data sample  $\mathbf{x}_0$  gradually loses its distinctive features as  $t$  increases. When  $T$  approaches infinity,  $\mathbf{x}_T$  converges to an isotropic Gaussian distribution.

An advantage of this process is that we can obtain samples at any arbitrary time step using a closed-form expression with the reparameterization trick.

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \bar{\boldsymbol{\epsilon}}_{t-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon} \end{aligned}$$

We can come to the following equation:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (2)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ ,  $\boldsymbol{\epsilon}_{t-1}, \boldsymbol{\epsilon}_{t-2}, \dots \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\bar{\boldsymbol{\epsilon}}_{t-2}$  merges two Gaussians.

**Reverse Diffusion Process** Reversing this process and sampling from  $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$  would enable us to reconstruct the true sample from a Gaussian noise input,  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . If  $\beta_t$  is sufficiently small,  $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$  will also be Gaussian. However, estimating  $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$  is not feasible because it requires using the entire dataset. Therefore, we need to learn a model  $p_\theta$  that approximates these conditional probabilities for running the reverse diffusion process.

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)) \quad (3)$$

Where  $\theta$  is a learnable parameter vector in the Gaussian distribution's mean function  $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$  and standard deviation function  $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)$ , the data samples generated by this distribution can be represented as:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) \quad (4)$$

By applying the learned parameters  $\theta$  to the mean function  $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$  and the standard deviation function  $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)$ . In brief, the forward diffusion process adds noise to the data sample, while the reverse diffusion process removes the noise and creates new data samples.

**Training Objective of Diffusion Probabilistic Models** Similar to Variational Autoencoder [35], the variational lower bound can be used to optimize the negative log-likelihood as follows:

$$\begin{aligned} -\log p_\theta(\mathbf{x}_0) &\leq -\log p_\theta(\mathbf{x}_0) + D_{\text{KL}}(q(\mathbf{x}_{1:T} | \mathbf{x}_0) \| p_\theta(\mathbf{x}_{1:T} | \mathbf{x}_0)) \\ &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_{\mathbf{x}_{1:T} \sim q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T}) / p_\theta(\mathbf{x}_0)} \right] \\ &= -\log p_\theta(\mathbf{x}_0) + \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} + \log p_\theta(\mathbf{x}_0) \right] \\ &= \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \end{aligned}$$

We can come to the following equations:

$$L_{\text{VLB}} = \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \geq -\mathbb{E}_{q(\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0) \quad (5)$$

To make each term in the equation analytically computable, we can reformulate the objective as a combination of several terms involving KL divergence and entropy. The objective can be rewritten as follows:

$$\begin{aligned}
 L_{\text{VLB}} &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[ \log \frac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \\
 &= \mathbb{E}_q \left[ \log \frac{\prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1})}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} \right] \\
 &= \mathbb{E}_q \left[ -\log p_\theta(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} + \log \frac{q(\mathbf{x}_T | \mathbf{x}_0)}{q(\mathbf{x}_1 | \mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1 | \mathbf{x}_0)}{p_\theta(\mathbf{x}_0 | \mathbf{x}_1)} \right] \\
 &= \mathbb{E}_q \left[ \log \frac{q(\mathbf{x}_T | \mathbf{x}_0)}{p_\theta(\mathbf{x}_T)} + \sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)}{p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)} - \log p_\theta(\mathbf{x}_0 | \mathbf{x}_1) \right] \\
 &= \mathbb{E}_q \left[ \underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} \right. \\
 &\quad \left. - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0} \right]
 \end{aligned}$$

Since  $x_0$  follows a fixed data distribution and  $x_T$  is a Gaussian noise,  $L_T$  is a constant term. We can also interpret  $L_0$  as the entropy of the multivariate Gaussian distribution, because  $p_\theta(x_0 | x_1)$  is a Gaussian distribution with mean  $\mu_\theta(x_1, 1)$  and covariance matrix  $\Sigma_\theta$ . The loss term  $L_t, t \in [1, 2, 3, \dots, T-1]$  can be parameterized as:

$$L_t = E_{x_0, \epsilon_t} \left[ \frac{\beta_t^2}{2\alpha_t(1-\bar{\alpha}_t)\Sigma_\theta^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon_t, t)\|^2 \right] + C \quad (6)$$

According to [24], the diffusion model can be trained more effectively by using a simplified objective that does not include the weighting term:

$$L_{\text{simple}}(\theta) = E_{x_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon_t, t)\|^2 \right] + C \quad (7)$$

**Classifier-Free Guidance for Conditional Generation** Conditional diffusion steps can be performed by combining the scores from both a conditional and an unconditional diffusion model [36]. The unconditional diffusion probabilistic model  $p_\theta(\mathbf{x})$  is parameterized by a score estimator  $\epsilon_\theta(\mathbf{x}_t, t)$ , while the conditional model  $p_\theta(\mathbf{x} | y)$  is parameterized by  $\epsilon_\theta(\mathbf{x}_t, t, y)$ . A single neural network can learn these two models simultaneously.

The implicit classifier’s gradient can be expressed with conditional and unconditional score estimators. The classifier-guided modified score, which incorporates this gradient, does not depend on a separate classifier.

$$\nabla_{\mathbf{x}_t} \log p(y | \mathbf{x}_t) = -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} (\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, y) - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)) \quad (8)$$

$$\bar{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, t, y) = (w + 1)\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, y) - w\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \quad (9)$$

In this study, we present the first study on applying the classifier-free guidance diffusion probabilistic model to generate high-quality data samples for environmental sound classification tasks.

## 2.2 DPM-Solver++ and DPM-Solver

One of the main challenges of working with diffusion probabilistic models is the high computational cost and time required to generate data samples from the complex posterior distributions. To overcome this limitation, we adopt DPM-Solver++ as our sampling method.

DPM-Solver is a high-order solver that can generate high-quality samples in around 10 steps by solving the diffusion ODE with a data prediction model. DPM-Solver++ is an improved version of DPM-Solver that can handle guided sampling by using thresholding methods to keep the solution matching the training data distribution.

The DPM-Solver++ and DPM-Solver raise the efficiency of training-free samplers to a new level to generate high-quality samples in the "few-step sampling" regime. This is the regime in which the sampling can be done within approximately 10 steps of sequential function evaluations. The DPM-Solver++ and DPM-Solver tackle the alternative problem of sampling from DPMs by solving the corresponding diffusion ordinary differential equations (ODEs) of DPMs. Moreover, diffusion ODEs have a semi-linear structure, meaning they comprise two parts: a linear function dependent on the data variable and a nonlinear function parameterized by neural networks. Consequently, the DPM-Solver++ and DPM-Solver use a precise formulation of the solutions of diffusion ODEs by analytically computing the linear portion of the solutions, thereby preventing the discretization mistake that the corresponding discretization would otherwise cause. In addition, it is possible to simplify the solutions to an exponentially weighted summation of the neural network by employing change-of-variable. This can be done efficiently. Such an integral is very special and can be efficiently approximated by numerical methods for exponential integrators.

Customized solver for diffusion ODEs is shown below:

$$\mathbf{x}_{t_{i-1} \rightarrow t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \sum_{n=0}^{k-1} \hat{\boldsymbol{\epsilon}}_\theta^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}}) \int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda \quad (10)$$

Where  $\mathbf{x}_s$  is an initial value at time  $s > 0$ ,  $\mathbf{x}_t$  is the solution at time  $t \in [0, s]$  and  $\lambda_t := \log(\alpha_t / \Sigma_t)$ . As  $\lambda(t) = \lambda_t$  is a strictly decreasing function of  $t$ , it has an inverse function  $t_\lambda(\cdot)$  satisfying  $t = t_\lambda(\lambda(t))$ . The DPM-Solver++ further changes the subscripts of  $\mathbf{x}$  and  $\boldsymbol{\epsilon}_\theta$  from  $t$  to  $\lambda$  and denote

$\hat{\mathbf{x}}_\lambda := \mathbf{x}_{t_\lambda(\lambda)}$ ,  $\hat{\epsilon}_\theta(\hat{\mathbf{x}}_\lambda, \lambda) := \epsilon_\theta(\mathbf{x}_{t_\lambda(\lambda)}, t_\lambda(\lambda))$ . The  $\mathcal{O}(h_i^{k+1})$  is omitted in the above equation because it is a high-order error.

The main difference between DPM-Solver and DPM-Solver++ is that the former is designed for sampling without guidance, while the latter is designed for sampling with guidance. Sampling with guidance means using additional information, such as text or images, to guide the sampling process of DPMs, which can improve the sample quality and diversity.

The primary technique that DPM-Solver++ uses to adapt to guided sampling is thresholding. Thresholding is a method to keep the solution of the diffusion ODE within the training data distribution by applying a hard or soft threshold to the pixel values or the latent variables. Thresholding can reduce the noise and ambiguity in the generated images and improve the sample quality and diversity. DPM-Solver++ adopts two types of thresholding: dynamic thresholding and static thresholding. Dynamic thresholding adjusts the threshold value according to the noise level and the guidance scale, while static thresholding uses a fixed threshold value for all steps. DPM-Solver++ combines both types of thresholding to balance the trade-off between stability and efficiency.

### 2.3 Data Augmentation

Data augmentation is a potent tactic to broaden the current data range and enable model training without requiring new data collection. In this research, standard and intelligent data augmentation methodologies are also taken into consideration. Two distinct audio data deformations are applied in conventional data augmentation.

First, certain background noises were added to the data samples, including crowd, street, and restaurant sounds (the background noises were taken from publicly available recordings made available on the "freesound.org" website [25]). Second, the records are subjected to pitch shifting [26]. The audio samples' pitch is adjusted by a half-octave (up and down) to produce various sounds. The audio stream is subjected to each contortion before being transformed into the input representation. Thirdly, the audio sample augmentation implements the time stretch [30].

For intelligent data augmentation, we use the U-net structure in the [28] for diffusion probabilistic models and DPM-Solver++ for the sampling schedule.

### 2.4 Top-k Selection Pretrained Discriminator

One of the challenges of using DPMs for data augmentation is that the quality of the generated samples may vary depending on the amount of available data and computational resources. To address this issue, we propose to use a pretrained discriminator network to filter out the low-quality samples and retain only the ones that are realistic and diverse enough to augment the training data. The discriminator network is an Xception [37] trained on the entire dataset to classify the images into their corresponding labels. The filtering criterion is based on the top-k accuracy of the discriminator, i.e., we accept a generated sample if its

accurate label is among the top-k predictions of the discriminator. Otherwise, we reject it. This way, we ensure that the generated samples are visually plausible and semantically consistent with their labels. The number of accepted samples can be expressed as follows:

$$G = \sum_{i=1}^N \mathbb{I}(f_k(\mathbf{x}_i, c_i) = c_i) \quad (11)$$

where  $c_i$  is the label of the  $i$ -th generated sample,  $\mathbf{x}_i$  is the  $i$ -th generated sample,  $f_k$  is the discriminator network with top-k prediction, and  $N$  is the number of generation epochs.

## 2.5 DL models for Environmental Sound Classification

One of the main objectives of this study is to assess the quality of the synthetic data samples generated by DPMs for environmental classification tasks. To this end, we propose to use the synthetic data samples to augment the original training data, and then train different deep learning (DL) classifiers on the augmented data. We hypothesize that the augmented data can enhance the diversity and robustness of the training data, and thus improve the performance of the DL classifiers for weed recognition. To test this hypothesis, we select seven state-of-the-art DL models from different architectures and paradigms, namely ResNet-50 [27], Xception [37], ConViT-tiny [38], mobilevitv2-50 [39], mobilevitv2-150 [39], ConvNext-tiny [40] and Deit III [41]. These models are implemented using the timm library [31], and their hyperparameters are set to their default values as suggested by the authors. We evaluate the performance of these models on the UrbanSound8K dataset.

# 3 Experiments

## 3.1 Experiments Pipeline

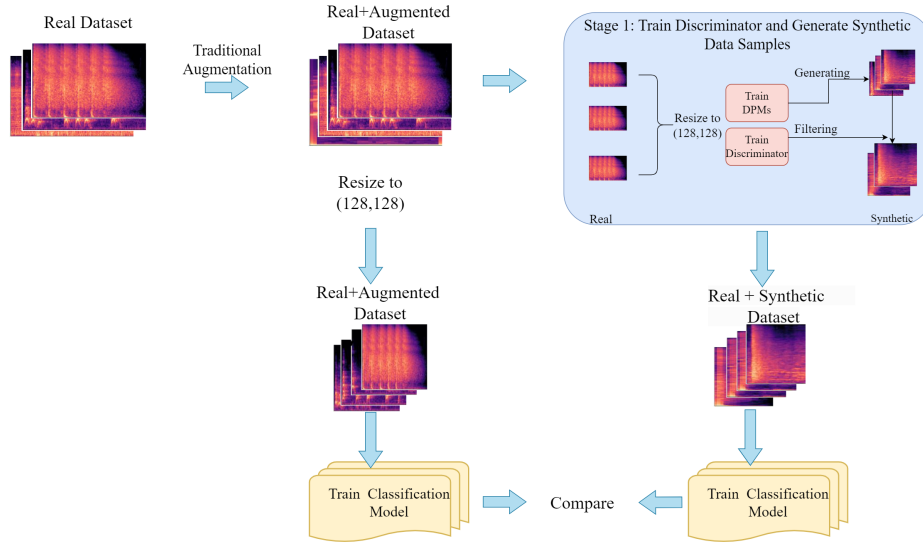
The whole Pipeline is shown in Fig 1.

The procedure can be divided into several steps in the proposed experimental settings. Initially, the audio is transformed into mel-spectrogram, which serves as the foundation for the subsequent steps. Following this, the training samples are augmented and resized into  $128 \times 128$ , and the resulting samples are incorporated into the original dataset. Simultaneously, a discriminator model is trained on the augmented dataset.

Once the augmented dataset is prepared, a diffusion probabilistic model is trained on it. This model is then employed to generate data samples. The discriminator model, which was previously trained, is used to filter out unqualified data samples from the generated samples. Finally, **8730 qualified synthetic data samples** are generated.

Lastly, two comparison experiments are conducted to evaluate the effectiveness of the models. The first experiment involves training a model using the





**Fig. 1.** An overview of the proposed pipeline for generating synthetic data samples. The pipeline consists of four main components: data augmentation, diffusion probabilistic modeling, discriminator filtering, and comparison experiments.

original data samples and traditional augmentation techniques. In contrast, the second experiment trains a model using the original data samples along with synthetic data samples.

### 3.2 UrbanSound8K Dataset

UrbanSound8K is an audio dataset that contains 8732 labelled sound excerpts ( $\leq 4s$ ) of urban sounds from 10 classes: air conditioner, car horn, children playing, dog bark, drilling, engine idling, gunshot, jackhammer, siren, and street music. The classes are drawn from the urban sound taxonomy. All excerpts are taken from field recordings uploaded to [www.freesound.org](http://www.freesound.org).

The dataset can be used for various tasks such as urban sound classification, sound event detection, acoustic scene analysis, etc.

### 3.3 Hyperparameters Setting

**Data Preprocessing** This section investigates the performance of the diffusion probabilistic model using the UrbanSound8K dataset to augment the data samples and optimize the classification model for environmental sound recognition. The data samples have a length of up to 5 seconds. The MFCC feature extraction is utilized, and the mel-spectrograms are generated. The original images have a dimension of  $768 \times 384$ . To minimize the number of training parameters in the diffusion probabilistic model training process, the original images are resized to  $128 \times 128$ , where each image has 128 frames (columns) and 128 bands (rows).

**Hyperparameters Setting for Data Augmentation** To enhance the stability of model training, we apply each data augmentation technique to the data samples in a stochastic manner described by the hyperparameter called  $p$ . Each data sample will undergo no more than two random transformations and at least one transformation. The specific setting is presented in Table 1.

The city ambience noise is composed of three pieces of sound. Each sound has a duration of 60 seconds. To augment the data sample with this noise, we randomly select a segment of the noise that has the same length as the data sample. Then, we superimpose this segment on the original data sample by adding their amplitudes. This creates a new data sample that contains both the original signal and the city ambience noise.

Pitch shifting refers to altering a sound’s pitch by increasing or decreasing its frequency. The pitch shift factor can quantify the degree of pitch shifting, which is defined as the ratio between the output and input frequencies. For instance, a pitch shift factor of 2 indicates that the output frequency is double the input frequency. Pitch shifting can be performed in two directions: up pitch shifting, which increases the frequency and raises the pitch; and down pitch shifting, which decreases the frequency and lowers the pitch. The transformation of up and down pitch shifting will not be implemented simultaneously.

Time stretch is changing the speed or duration of an audio signal without affecting its pitch. The minimum rate and maximum rate parameters control the minimum and maximum speed change factor, respectively. For example, a minimum rate of 0.8 means 20% can slow the audio down, and a maximum rate of 1.25 means the audio can be sped up by 25%.

**Table 1.** Parameters Setting for Traditional Data Augmentation

Method	$p$	Setting
Noise from City Ambience	0.6	The weight of the ambience is 0.6.
Up Pitch Shifting	0.8	The pitch shift factor is 2.
Down Pitch Shifting	0.8	The pitch shift factor is 2.
Time stretch	0.7	Minimum rate is 0.8 and Maximum rate is 1.25

**Diffusion Probabilistic Model Training Setting** In this paper, we use the U-net structure implemented in the [28] for the estimation of  $\mathbf{p}_\theta$ . U-net performs four upsampling operations in its decoder path and uses skip connections to concatenate feature maps from the same stage in the encoder path instead of directly supervising and backpropagating loss on high-level semantic features[29]. This ensures that the final recovered feature map integrates more low-level features that capture fine details and enables the fusion of features at different scales for multi-scale prediction. The four upsampling operations also make the image

recover edge information more finely by reducing spatial resolution loss. Consequently, the U-net structure suffers much less information loss compared with other methods when sampling. As a result, the U-net structure is preferred in the diffusion probabilistic model.

The Unet architecture was configured with the following parameters: the number of channels in the first convolutional layer was set to 64 (dim=64), the number of channels was multiplied by 1, 2, 4, and 8 in the subsequent down-sampling and upsampling layers (dim\_mults=(1,2,4,8)), the number of residual blocks in each group was fixed at 8 (resnet\_block\_groups=8), and the dimensionality of the learned sinusoidal embeddings was chosen as 16 (learned\_sinusoidal\_dim=16).

The optimization algorithm used for training the U-Net network is AdamW [33], a variant of Adam [34] incorporating weight decay regularization. The learning rate is set to 0.0001, a standard value for deep learning models. The weight decay parameter is set to 0.05, which helps to prevent overfitting by penalizing large weights. The other parameters of AdamW, such as beta values and epsilon values, are kept at their default values as suggested by the original paper [33].

The loss function used to measure the discrepancy between the predicted and ground truth images is the mean squared error (MSE), defined as the average of the squared differences between the pixel values. The MSE is a widely used loss function for image reconstruction tasks, as it encourages high-fidelity reconstructions.

The number of training epochs is set to 3500 for the dataset with augmented data, which is sufficient for the network to converge to a stable solution.

To demonstrate the effectiveness of our network, we refrain from using techniques such as model transfer, exponential moving average (EMA), pretraining, or other tricks.

**Diffusion Probabilistic Model Sampling Setting** There are two versions of DPM-Solver++: one is DPM-Solver++ (2S), which is a second-order single-step solver, and the other is DPM-Solver++ (2M), which is a multistep second-order solver. The latter deals with the instability problem of the high-order solvers by reducing the effective step size. In our paper, we use the 2M as the sampling schedule in our study. The method is implemented via diffusers[42].

The following settings are used for the DPM-Solver++ parameters:

The initial and final values of  $\beta$  for inference are 0.0001 and 0.02, respectively.  $\beta$  is a hyperparameter that regulates the balance between the data likelihood and the prior distribution over the latent variables. The latent variables are unobserved variables that capture the underlying structure of the data.

We use the linear method for the  $\beta$  schedule that maps a range of betas to a sequence of betas for updating the model. The linear method progressively increases  $\beta$  from the initial value to the final value over a predetermined number of iterations. The order of the DPM-Solver++ is 2, which indicates that it employs a second-order differential equation to model the dynamics of the latent variables.

The solver type for the second-order solver is the midpoint, which is a numerical method that approximates the solution of the differential equation by using the midpoint of an interval as an estimate of its values.

The number of inference steps is 20. The parameter means the number of diffusion steps used when generating samples with a pre-trained model.

Other hyperparameters are set as default.

**Classification Models’ Training Setting** We trained our model for 500 epochs with a batch size of 30. We used the AdamW optimizer with the same hyperparameters as the DPMs training optimizer, such as learning rate, weight decay and epsilon. We used the cross-entropy loss function with label smoothing of 0.1 to prevent overfitting and improve generalization. We kept the other hyperparameters as default in the timm, such as dropout rate, hidden size, the number of layers, etc. The hyperparameter  $k$  of top-k selection is set as one in the following experiments. The experiments are performed on a computer with a 13th Gen Intel R Core™ i9-13900KF CPU and a GeForce RTX 4090 GPU (24 GB GDDR6X memory).

**Discriminator Training Setting** This paper employs the Xception model implemented by timm [31] as the backbone of the proposed method. The Xception model is a deep convolutional neural network that can achieve high accuracy on various image recognition tasks. The following parameters are used to configure the Xception model for this paper:

We used the Xception model with the following hyperparameters: number of input channels = 3, dropout rate = 0., and global pooling method = average pooling.

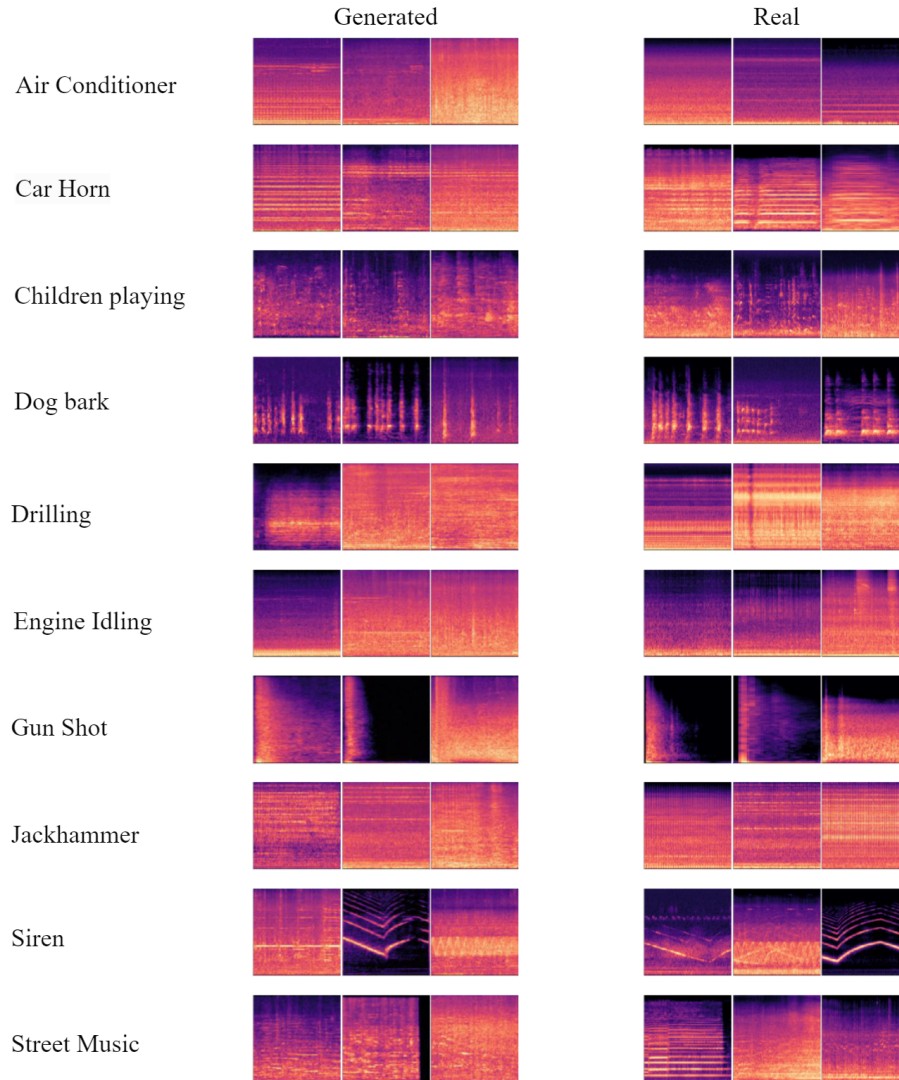
We used the same training hyperparameters as the classification models, such as learning rate, batch size, number of epochs and optimizer.

### 3.4 Experiments Results

**Random Visual Samples of Generated Data** Fig. 2 illustrates random visual examples of the generated spectrograms using DPMs. As you can see in this figure, DPMs have a high capability to produce spectrograms that have similar structures.

**Results of Augmentation for Different SOTA Models** To assess the testing performance of the DL models and their generalization ability on unseen data, we use the 10-fold method as a fair and reliable evaluation technique. The final accuracy is computed as the mean of the 10 test results.

Table 2 shows the testing performance of seven DL models, ResNet-50 [27], Xception [37], ConViT-tiny [38], Mobilevitv2-50 [39], Mobilevitv2-150 [39], Conv Next-tiny [40] and Deit III [41], for environmental sound classification on UrbanSound8K dataset with synthetic augmentation + real dataset and with traditional data augmentation + real dataset.



**Fig. 2.** Real (right panel) and generated (left panel) audio samples intelligent augmentation. Each row represents one sound class.

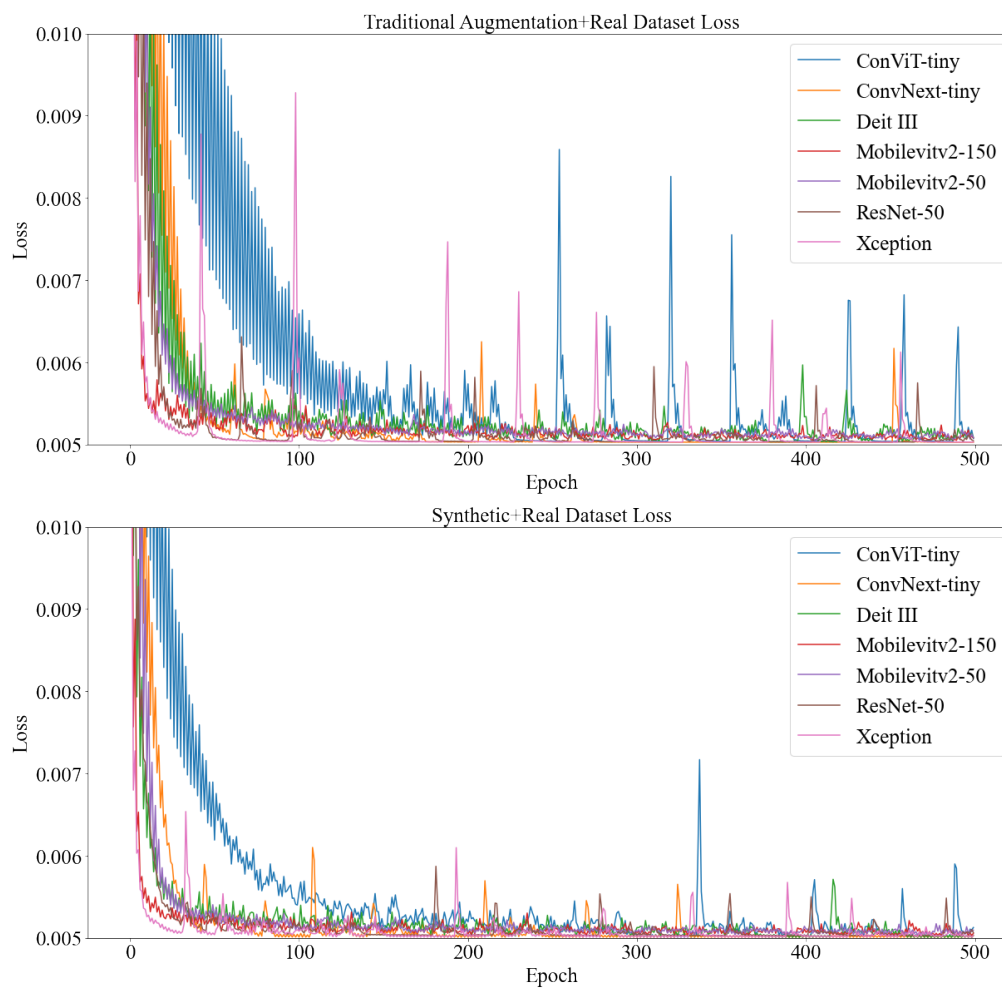
It is evident that considerable improvements have been achieved for all DL models by incorporating expanded datasets with synthetic images, in terms of more stable training (see loss curves in Fig 3) and higher accuracy and lower losses. For example, with the synthetic augmented images, Inception-v3 and ResNet-50 yield classification accuracies of 80.1% and 80.2% on the UrbanSound8K dataset, respectively representing about 6.3% and 7.6% improvements over the baseline models without synthetic augmentation.

**Table 2.** Performance comparison of DL models on environmental sound classification trained with and without the samples generated with the proposed data augmentation. The best values are in bold.

Index	Models	Parameters	UrbanSound8K
			Top-1 Accuracy (%)
Real+Traditional	ResNet-50	23528522	73.8%
	Xception	20827442	72.5%
	ConViT-tiny	5494098	65.1%
	Mobilevitv2-50	1116163	59.2%
	Mobilevitv2-150	9833443	68.3%
	ConvNext-tiny	27827818	60.7%
	Deit III	85722634	67.4%
Real+Synthetic	ResNet-50	23528522	<b>80.1%</b>
	Xception	20827442	<b>80.2%</b>
	ConViT-tiny	5494098	<b>68.6%</b>
	Mobilevitv2-50	1116163	<b>62.0%</b>
	Mobilevitv2-150	9833443	<b>74.6%</b>
	ConvNext-tiny	27827818	<b>65.9%</b>
	Deit III	85722634	<b>72.2%</b>

**Influence of Hyperparameter  $k$  for Top-k Selection** To assess how the hyperparameter  $k$  affects the outcome of synthetic data generation, we perform a series of experiments with different values of  $k$  from 1 to 10 and compare the results using accuracy. The results are shown in Table 4.

The results show that the top-k selection strategy effectively enhances most DL models' performance. By selecting the most confident synthetic images for each class, the top-k strategy reduces the noise and ambiguity in the augmented data. As shown in Table 4, the top-k strategy significantly improves accuracy for six models. For instance, ResNet-50 achieves an accuracy of 80.1% with  $k=1$ , which is 2.7% higher than the model without top-k selection. Similarly, Xception



**Fig. 3.** Training loss curves of different DL models on the UrbanSound8K dataset with and without synthetic data augmentation. The left figure represents the baseline models trained without synthetic augmentation, while the right figure represents those trained with synthetic augmentation. The results show that synthetic augmentation helps to reduce the loss and improve the stability of the training process for all models.

**Table 3.** The impact of hyperparameter  $k$  on the accuracy of seven deep learning models on a synthetic+real dataset. The table shows how the accuracy changes as  $k$  varies from 1 to 10. The value of  $k$  does not affect the generation of data samples when  $k \geq 5$  in our training setting. Therefore, the accuracy of the models is constant for these values of  $k$ . The best values are in bold.

Method	Different Values of Hyperparameter $k$				
	1	2	3	4	$\geq 5$
ResNet-50	<b>80.1%</b>	80.2%	77.3%	76.7%	77.4%
Xception	<b>80.2%</b>	78.7%	74.8%	74.5%	74.0%
ConViT-tiny	68.6%	<b>68.7%</b>	64.5%	64.5%	64.6%
Mobilevitv2-50	62.0%	62.8%	63.0%	65.5%	<b>69.0%</b>
Mobilevitv2-150	74.6%	73.5%	73.5%	<b>76.4%</b>	72.6%
ConvNext-tiny	<b>65.9%</b>	62.3%	62.0%	61.0%	62.8%
Deit III	72.2%	<b>72.8%</b>	71.3%	69.2%	67.8%

attains an accuracy of 80.2% with  $k=1$ , which is 6.2% higher than the models without top- $k$  selection.

## 4 Conclusion

This paper introduced a novel application of diffusion models for generating high-quality synthetic images from environmental sound recordings. To the best of our knowledge, this is the first study that explores the use of diffusion models for data augmentation in environmental sound classification. We also proposed a new selection method based on the top- $k$  confidence scores to filter out the low-quality synthetic images and retain the most informative ones for each sound class.

We conducted extensive experiments on a widely used sound dataset, Urban-Sound8K, and evaluated the performance of seven state-of-the-art DL models trained on the augmented datasets with different settings. The experimental results demonstrated that the diffusion models can generate realistic and diverse synthetic images that can effectively improve the classification accuracy and reduce the losses for all DL models. Moreover, the top- $k$  selection method further enhanced the performance by removing noisy and ambiguous synthetic images and increasing the data balance among classes.



## References

1. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33, 6840–6851 (2020)
2. Salamon, Justin and Juan Pablo Bello. “Deep Convolutional Neural Networks and Data Augmentation for Environmental Sound Classification.” *IEEE Signal Processing Letters* 24 (2016): 279-283.
3. Gong, Yuan, Yu-An Chung and James R. Glass. “AST: Audio Spectrogram Transformer.” *ArXiv abs/2104.01778* (2021): n. pag.
4. Chen, Ke, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick and Shlomo Dubnov. “HTS-AT: A Hierarchical Token-Semantic Audio Transformer for Sound Classification and Detection.” *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2022): 646-650.
5. Bahmei, Behnaz, Elina Birmingham and Siamak Arzanpour. “CNN-RNN and Data Augmentation Using Deep Convolutional Generative Adversarial Network for Environmental Sound Classification.” *IEEE Signal Processing Letters* 29 (2022): 682-686.
6. Hershey, Shawn, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, R. Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron J. Weiss and Kevin W. Wilson. “CNN architectures for large-scale audio classification.” *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2016): 131-135.
7. Shrestha, Roman, Cornelius Glackin, Julie A. Wall and Nigel Cannings. “Bird Audio Diarization with Faster R-CNN.” *International Conference on Artificial Neural Networks* (2021).
8. X. Zhu, Y. Liu, J. Li, T. Wan, and Z. Qin, “Emotion classification with data augmentation using generative adversarial networks,” in *Proc. PacificAsia Conf. Knowl. Discov. Data Mining*, Cham, Switzerland: Springer, Jun. 2018, pp. 349–360, doi: 10.1007/978-3-319-93040-4\_28.
9. Frid-Adar, Maayan, Eyal Klang, Michal Marianne Amitai, Jacob Goldberger and Hayit Greenspan. “Synthetic data augmentation using GAN for improved liver lesion classification.” *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)* (2018): 289-293.
10. Lee, Han S., Haeil Lee, Helen Hong, Heejin Bae, Joon Seok Lim and Junmo Kim. “Classification of focal liver lesions in CT images using convolutional neural networks with lesion information augmented patches and synthetic data augmentation.” *Medical physics* (2021): n. pag.
11. Arjovsky, Martín, Soumith Chintala and Léon Bottou. “Wasserstein GAN.” *ArXiv abs/1701.07875* (2017): n. pag.
12. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1):53–65, 2018.
13. Mescheder, S. Nowozin, and A. Geiger. The numerics of gans. *Advances in neural information processing systems*, 30, 2017.
14. Zhao, H. Ren, A. Yuan, J. Song, N. Goodman, and S. Ermon. Bias and generalization in deep generative models: An empirical study. *Advances in Neural Information Processing Systems*, 31, 2018.
15. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

16. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
17. . Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780– 8794, 2021.
18. Müller-Franzes, Gustav, Jan Moritz Niehues, Firas Khader, Soroosh Tayebi Arasteh, Christoph Haarbuerger, Christiane Kuhl, Tian Wang, Tianyu Han, Sven Nebelung, Jakob Nikolas Kather and Daniel Truhn. “Diffusion Probabilistic Models beat GANs on Medical Images.” *ArXiv abs/2212.07501* (2022): n. pag.
19. Maz’e, Francois and Faez Ahmed. “Diffusion Models Beat GANs on Topology Optimization.” (2022).
20. Song, Jiaming, Chenlin Meng and Stefano Ermon. “Denoising Diffusion Implicit Models.” *ArXiv abs/2010.02502* (2020): n. pag.
21. Lu, Cheng, Zhou, Yuhao, Bao, Fan, Chen, Jianfei, Li, Chongxuan, and Jun Zhu. "DPM-Solver++: Fast Solver for Guided Sampling of Diffusion Probabilistic Models." *ArXiv*, (2022). Accessed March 22, 2023. /abs/2211.01095.
22. Cordts, Marius, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth and Bernt Schiele. “The Cityscapes Dataset for Semantic Urban Scene Understanding.” *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016): 3213-3223.
23. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
24. Saharia, Chitwan, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet and Mohammad Norouzi. “Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding.” *ArXiv abs/2205.11487* (2022): n. pag.
25. F. Font, G. Roma, and X. Serra, “Freesound technical demo,” in *Proc. 2013 ACM Multimedia Conf., Spain, 2013*, pp. 411–412, doi: 10.1145/2502081.2502245
26. J. Salamon and J. P. Bello, “Deep convolutional neural networks and data augmentation for environmental sound classification,” *IEEE Signal Process. Lett.*, vol. 24, no. 3, pp. 279–283, Mar. 2017, doi: 10.1109/LSP. 2017.2657381.
27. K. He, X. Zhang, S. Ren and J. Sun, "Deep Residual Learning for Image Recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016, pp. 770-778, doi: 10.1109/CVPR.2016.90.
28. lucidrains.2023.lucidrains.denoising-diffusion-pytorch.  
[https://github.com/lucidrains/denoising-diffusion-pytorch\(2023\)](https://github.com/lucidrains/denoising-diffusion-pytorch(2023)).
29. Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional Networks for Biomedical Image Segmentation. *ArXiv*.  
<https://doi.org/10.48550/arXiv.1505.04597>
30. Iwana, Brian Kenji and Seiichi Uchida. “An empirical survey of data augmentation for time series classification with neural networks.” *PLoS ONE* 16 (2020): n. pag.
31. rw2019timm, Ross Wightman, 2019, PyTorch Image Models <https://github.com/rwightman/pytorch-image-models>
32. Press, William H.; Teukolsky, Saul A.; Vetterling, William T.; Flannery, Brian P. (1992). *Numerical recipes in C: the art of scientific computing* (2nd ed.). New York, NY, USA: Cambridge University Press. pp. 123-128. ISBN 0-521-43108-5.
33. Loshchilov, Ilya and Frank Hutter. “Decoupled Weight Decay Regularization.” *International Conference on Learning Representations* (2017).
34. Kingma, Diederik P. and Jimmy Ba. “Adam: A Method for Stochastic Optimization.” *CoRR abs/1412.6980* (2014): n. pag.

35. Kingma, Diederik P., and Max Welling. "Auto-Encoding Variational Bayes." ArXiv, (2013). Accessed March 22, 2023. /abs/1312.6114.
36. Jonathan Ho & Tim Salimans. "Classifier-Free Diffusion Guidance." NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications.
37. Chollet, François. "Xception: Deep Learning with Depthwise Separable Convolutions." 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016): 1800-1807.
38. d'Ascoli, Stéphane, Hugo Touvron, Matthew L. Leavitt, Ari S. Morcos, Giulio Biroli and Levent Sagun. "ConViT: improving vision transformers with soft convolutional inductive biases." Journal of Statistical Mechanics: Theory and Experiment 2022 (2021): n. pag.
39. Mehta, Sachin, and Mohammad Rastegari. "MobileViT: Light-weight, General-purpose, and Mobile-friendly Vision Transformer." ArXiv, (2021). Accessed March 23, 2023. /abs/2110.02178.
40. Liu, Zhuang, Hanzi Mao, Chaozheng Wu, Christoph Feichtenhofer, Trevor Darrell and Saining Xie. "A ConvNet for the 2020s." 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022): 11966-11976.
41. Touvron, Hugo, Matthieu Cord and Hervé Jégou. "DeiT III: Revenge of the ViT." European Conference on Computer Vision (2022).
42. von-platen-et-al-2022-diffusers, Patrick von Platen and Suraj Patil and Anton Lozhkov and Pedro Cuenca and Nathan Lambert and Kashif Rasul and Mishig Davaadorj and Thomas Wolf, 2022, Diffusers: State-of-the-art diffusion models, <https://github.com/huggingface/diffusers>
43. Chen, Yunhao, Zhu, Yunjie, Yan, Zihui, Shen, Jianlu, Ren, Zhen, and Yifan Huang. "Data Augmentation for Environmental Sound Classification Using Diffusion Probabilistic Model with Top-k Selection Discriminator." ArXiv, (2023). Accessed March 28, 2023. /abs/2303.15161.