

UNSUPERVISED VOCAL DEREVERBERATION WITH DIFFUSION-BASED GENERATIVE MODELS

Koichi Saito Naoki Murata Toshimitsu Uesaka Chieh-Hsin Lai
Yuhta Takida Takao Fukui Yuki Mitsufuji

Sony Group Corporation, Tokyo, Japan

ABSTRACT

Removing reverb from reverberant music is a necessary technique to clean up audio for downstream music manipulations. Reverberation of music contains two categories, *natural reverb*, and *artificial reverb*. Artificial reverb has a wider diversity than natural reverb due to its various parameter setups and reverberation types. However, recent supervised dereverberation methods may fail because they rely on sufficiently diverse and numerous pairs of reverberant observations and retrieved data for training in order to be generalizable to unseen observations during inference. To resolve these problems, we propose an unsupervised method that can remove a general kind of artificial reverb for music without requiring pairs of data for training. The proposed method is based on diffusion models, where it initializes the unknown reverberation operator with a conventional signal processing technique and simultaneously refines the estimate with the help of diffusion models. We show through objective and perceptual evaluations that our method outperforms the current leading vocal dereverberation benchmarks.

Index Terms—vocal dereverberation, diffusion-based generative models, weighted prediction error

1. INTRODUCTION

Reverb is one of the major audio effects that enable listeners to perceive the spatial characteristics, timbre, and texture of music. Reverb for music contains not only “natural reverb”, which has been studied extensively in the context of speech dereverberation problems, but also “artificial reverb”, which is mostly exploited as an effect for music production [1, 2]. In contrast, removing reverberated components from reverberant (wet) music signals is also an important technique for music production. Audio engineers and content creators do not merely mix up each signal when creating remixed, mastered, or upmixed music materials from existing contents—they may also apply new kinds of audio effects such as equalization or reverb [2–4] to original unprocessed (dry) music signals. However, whenever dry signals are not available, the reverb must be removed from the processed music signals before the desired manipulations can be applied.

To remove reverb from music signals, we need to take both natural and artificial reverb into account. Most research in this vein has focused on supervised approaches applied in a data-driven manner that require the preparation of numerous pairs of wet observations and their corresponding dry signals for training. This tends to get challenging because artificial reverb may have a higher number of variations than natural reverb due to its different parameter setups. More precisely, the creation of artificial reverb may involve various types of reverberations (e.g., *plate reverb*, *spring reverb*) and multiple parameters (e.g., *pre delay*, *decay rate*) in a set of reverb

having the same RT60. Supervised methods generally work well when target samples follow a similar distribution as the training data, but if they deviate from that distribution, the performance may degrade. Therefore, to remove artificial reverb, common supervised approaches may not work well for various types of wet signals because the training set may not exhaust the comprehensive data pairs (as discussed in Section 4).

Indeed, music dereverberation can be formulated as solving a linear inverse problem with a linear degradation operator. Recently, Denoising Diffusion Restoration Model (DDRM) [5], an unsupervised linear inverse problem solver based on diffusion-based generative models (i.e., diffusion models) [6, 7], has shown its effectiveness across various image restoration tasks. However, DDRM assumes a full knowledge of the linear degradation operator. This is problematic as in practical music dereverberation problems, linear degradation operators (i.e., reverberation operators) are generally unknown. Hence, directly applying DDRM to the music dereverberation problem may lead to inaccurate retrieval of the dry signal.

To resolve the above problems while (1) avoiding the collection of a large amount of paired data, (2) handling various types of reverb for music production (including artificial reverb), and (3) handling uncertainty of the reverberation operators, we propose an unsupervised method for music dereverberation using diffusion models that is motivated by DDRM. Our proposed method contains two key components. First, we extend DDRM into a practical scenario with unknown reverberation operators, where weighted prediction error (WPE) [8] is used to initially estimate a reverberation operator. Second, we propose adaptively correcting the initial reverberation operator estimated by WPE to obtain a more accurate one with the help of the predicted dry signal from a diffusion model. Our method only needs dry signals for training, which circumvents the needs to prepare the diverse data pairs required by supervised methods. We demonstrate that our method effectively dereverbs various types of wet vocal and outperforms unsupervised and supervised benchmarks through comprehensive objective and subjective evaluations on a set of wet vocal test data. We refer to <https://koichi-saito-sony.github.io/unsupervised-vocal-dereverb/> for examples of the generated audio samples, and reproducing our experiments with the detailed settings of the training configurations.

2. RELATED WORK

2.1. Existing research for music dereverberation

The early unsupervised approaches for music dereverberation focused primarily on signal processing techniques [9–14]. Some of these methods utilized linear prediction [11–13], while others were based on the assumption of a source model that directly expresses the harmonic structure of music signals [9, 10].

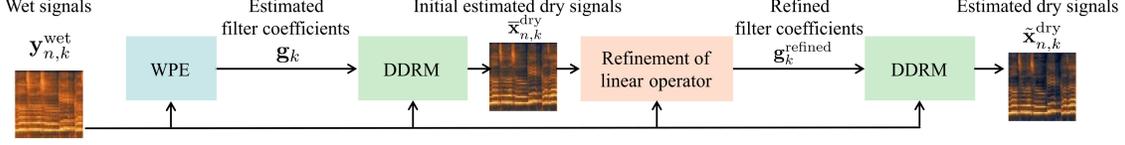


Fig. 1. Music dereverberation procedure of proposed method.

Some of the recent DNN-based approaches have achieved state-of-the-art performances on vocal and music dereverberation by using an end-to-end DNN [1] or by combining a conditional generative adversarial network with a diffusion-based vocoder [15]. Both methods must be trained with pairs of dry and wet signals and are shown to have a decent performance on their test datasets. However, due to the wide variety of artificial reverb discussed earlier, the effectiveness of these methods may degrade if they are applied to a test dataset that follows a different distribution of wet signals than the one the train dataset was derived from (discussed in more detail in Section 4).

2.2. Deep generative models for linear inverse problems

A wide class of image or audio restoration problems can be formulated as linear inverse problems [16–18]. Recent methods [5, 19–24] restore the information of data distribution via a deep generative model and use it as a prior to retrieve the signal from an observation. DDRM [5] aligns with this track in that a highly expressive pre-trained diffusion model is taken as the prior distribution, and DDRM exploits it as a general inverse problem solver with the assumption that the linear data degradation process is known. Since the training of diffusion models requires only clean data (dry signals), not pairs of clean and degraded data, only information about the degradation process is needed during inference.

3. PROPOSED METHOD

We propose an unsupervised music dereverberation method for unknown reverberation operators that exploits a pre-trained diffusion model on dry music signals as a prior and does not require pairs of wet and dry signals for training. The proposed method consists of three steps: (1) the initial estimation of the linear operator for the inverse problem via WPE [8] (Section 3.2), (2) the dry music estimation with DDRM (Section 3.3), and (3) the generative model-based refinement of the linear operator (Section 3.4).

3.1. Problem setting

Let $y_{n,k}^{\text{wet}} \in \mathbb{C}$ be wet signals in the short time Fourier transformation (STFT) domain, where n and k denote the time and frequency indices, respectively. We model the wet signals as

$$y_{n,k}^{\text{wet}} = x_{n,k}^{\text{dry}} + x_{n,k}^{\text{reverb}} + z_{n,k}, \quad (1)$$

where $x_{n,k}^{\text{dry}} \in \mathbb{C}$ and $x_{n,k}^{\text{reverb}} \in \mathbb{C}$ are the dry and reverb components included in the wet signal, respectively. Additive noise $z_{n,k} \in \mathbb{C}$ is assumed. The aim of music dereverberation is to estimate the dry signals $x_{n,k}^{\text{dry}}$ from the wet signals $y_{n,k}^{\text{wet}}$. Here, we assume that a generative model as a prior is available, namely, a pre-trained diffusion model trained on dry signals (see Section 3.3).

3.2. Estimation of linear operator in linear inverse problem

In this subsection, we revisit WPE [8] and introduce its interpretation as a linear inverse problem. The idea of WPE is to esti-

mate the reverb components and subtract them from the wet signals. It models the reverb components as the convolution of the filter $\mathbf{g}_k = [g_{1,k}, \dots, g_{L,k}]^T \in \mathbb{C}^L$ with the length of L and the wet signals, where $(\cdot)^T$ denotes the transpose of a matrix (or a vector), as

$$\hat{x}_{n,k}^{\text{dry}} = y_{n,k}^{\text{wet}} - \sum_{l=1}^L g_{l,k}^* y_{n-D-l+1,k}^{\text{wet}}, \quad (2)$$

where $\hat{x}_{n,k}^{\text{dry}}$ is the estimate of dry signals, $(\cdot)^*$ denotes the complex conjugate, and D is the prediction delay. The filter is obtained on the basis of various assumptions regarding dry signals, e.g., dry signals at time indices n and n' are assumed to be mutually uncorrelated when $|n - n'| > \delta$ for a certain constant $\delta > 0$. For more details, refer to the original paper [8]. This filter \mathbf{g}_k will be refined adaptively with the help of a diffusion model (see Section 3.4). With the obtained filter, the estimate of the dry signals is obtained by rewriting Eq. (2) as follows:

$$\hat{\mathbf{x}}_{n,k}^{\text{dry}} = (\tilde{\mathbf{I}} - \mathbf{G}_k) \mathbf{y}_{n,k}^{\text{wet}}, \quad (3)$$

where $\mathbf{y}_{n,k}^{\text{wet}} = [y_{n,k}^{\text{wet}}, \dots, y_{n-D-L-m+1,k}^{\text{wet}}]^T$, $\hat{\mathbf{x}}_{n,k}^{\text{dry}} = [\hat{x}_{n,k}^{\text{dry}}, \dots, \hat{x}_{n-m,k}^{\text{dry}}]^T$, and m is the number of processed samples. $\tilde{\mathbf{I}}$ and \mathbf{G}_k are Toeplitz matrices defined as

$$\tilde{\mathbf{I}} = [\mathbf{I}_{m \times m}, \mathbf{0}_{m \times (D+L-1)}], \quad (4)$$

$$\mathbf{G}_k = \begin{bmatrix} & g_{1,k}^* & g_{2,k}^* & \cdots & g_{L,k}^* & 0 & \cdots & 0 \\ \mathbf{0}_{m \times D} & 0 & g_{1,k}^* & g_{2,k}^* & \cdots & g_{L,k}^* & \ddots & \vdots \\ & \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & 0 \\ & 0 & \cdots & 0 & g_{1,k}^* & g_{2,k}^* & \cdots & g_{L,k}^* \end{bmatrix}, \quad (5)$$

where $\mathbf{I}_{m' \times m'}$ denotes the identity matrix of size m' , and $\mathbf{0}_{m' \times n'}$ is the $m' \times n'$ rectangular matrix with all the elements 0.

Now, the estimated dry signals $\hat{\mathbf{x}}_{n,k}^{\text{dry}}$ can be interpreted as the least-squares solution of the following linear inverse problem:

$$\mathbf{y}_{n,k}^{\text{wet}} = (\tilde{\mathbf{I}} - \mathbf{G}_k)^\dagger \hat{\mathbf{x}}_{n,k}^{\text{dry}} + \mathbf{z}_{n,k}, \quad (6)$$

where $\mathbf{z}_{n,k} = [z_{n,k}, \dots, z_{n-D-L-m+1,k}]^T$. \mathbf{A}^\dagger is the pseudo-inverse of a matrix \mathbf{A} and is defined as $\mathbf{A}^\dagger = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1}$, where $(\cdot)^H$ denotes the Hermitian transpose of the matrix.

3.3. DDRM for solving linear inverse problems

Our approach for estimating dry signals is motivated by DDRM [5]. It solves linear inverse problems in an unsupervised way that requires the linear operator only during inference and hence, pairs of the wet and dry signals are not needed during training for supervision. Therefore, music dereverberation with artificial reverb may benefit from this mechanism when the reverberation operation is known.

We review how the prior diffusion models are obtained. First, we train a diffusion model on the dry signal dataset in the STFT

Algorithm 1 Proposed algorithm for music dereverberation

Input: Wet signal $\mathbf{y}_{n,k}^{\text{wet}}$,Pre-trained diffusion model $p_\theta(\mathbf{x}^{\text{dry}})$, Number of refinements of the linear operator N_{refine} , Step size α , and Regularization parameter λ **Output:** Estimated dry signal $\tilde{\mathbf{x}}_{n,k}^{\text{dry}}$ Estimate the filter coefficient \mathbf{g}_k with WPE [8]Estimate the dry signal $\tilde{\mathbf{x}}_{n,k}^{\text{dry}}$ with DDRM [5, Eq. (8)]**for** $i = 1$ to N_{refine} **do**Refine \mathbf{g}_k with Eq.(10)**end for**Estimate the dry signal $\tilde{\mathbf{x}}_{n,k}^{\text{dry}}$ with DDRM on the refined $\mathbf{g}_k^{\text{refined}}$

domain. We denote a set of $x_{n,k}^{\text{dry}}$ as \mathbf{x}^{dry} . Diffusion models [6, 7, 25] are generative models with a Markov chain $\mathbf{x}_T^{\text{dry}} \rightarrow \dots \rightarrow \mathbf{x}_t^{\text{dry}} \rightarrow \dots \rightarrow \mathbf{x}_0^{\text{dry}} = \mathbf{x}^{\text{dry}}$ represented by the following joint distribution:

$$p_\theta(\mathbf{x}_{0:T}^{\text{dry}}) = p_\theta^{(T)}(\mathbf{x}_T^{\text{dry}}) \prod_{t=0}^{T-1} p_\theta^{(t)}(\mathbf{x}_t^{\text{dry}} | \mathbf{x}_{t+1}^{\text{dry}}), \quad (7)$$

where only $\mathbf{x}_0^{\text{dry}}$ is used to generate samples for \mathbf{x}^{dry} . For the training, a fixed approximated posterior is introduced to evaluate an evidence lower bound (ELBO) on the maximum likelihood objective:

$$q(\mathbf{x}_{1:T}^{\text{dry}} | \mathbf{x}_0^{\text{dry}}) = q^{(T)}(\mathbf{x}_T^{\text{dry}} | \mathbf{x}_0^{\text{dry}}) \prod_{t=0}^{T-1} q^{(t)}(\mathbf{x}_t^{\text{dry}} | \mathbf{x}_{t+1}^{\text{dry}}, \mathbf{x}_0^{\text{dry}}), \quad (8)$$

and we adopt the following ELBO objective, which is induced from the Gaussian parameterization for p_θ and q :

$$\mathbb{E}_{t,q(\mathbf{x}_0^{\text{dry}}, \mathbf{x}_t^{\text{dry}})} \left[\gamma_t \|\mathbf{x}_0^{\text{dry}} - f_\theta^{(t)}(\mathbf{x}_t^{\text{dry}})\|_2^2 \right], \quad (9)$$

where $f_\theta^{(t)}$ is a neural network that characterizes p_θ and estimates a noiseless data from a noisy data $\mathbf{x}_t^{\text{dry}}$, and γ_t are positive weighting coefficients determined by q .

The inference of DDRM with a pre-trained diffusion model is based on [5, Eq. (8)]. In particular, DDRM requires the singular value decomposition (SVD) of the linear operator, as $(\tilde{\mathbf{I}} - \mathbf{G}_k)^\dagger = \mathbf{U}_k \Sigma_k \mathbf{V}_k^H$. In our case, we only require the SVD of $(\tilde{\mathbf{I}} - \mathbf{G}_k)$, since the SVD of the pseudo-inverse becomes $\mathbf{V} \Sigma^\dagger \mathbf{U}^H$ if the SVD of the original matrix is $\mathbf{U} \Sigma \mathbf{V}^H$. Using the pre-trained diffusion-based model $p_\theta(\mathbf{x}^{\text{dry}})$ and the SVD of the linear operator, DDRM generates samples $\tilde{\mathbf{x}}_{n,k}^{\text{dry}}$ [5, Eq. (8)] that are consistent with the linear inverse problem in Eq. (3).

3.4. Generative model-based refinement of linear operator

Since WPE estimates the filter coefficients \mathbf{g}_k on the basis of relatively simple assumptions, e.g., the correlation property of dry signals, it does not necessarily provide reasonable filter coefficients, which are required in DDRM for the linear operator. We therefore refine the filter coefficients utilizing the dry signal $\tilde{\mathbf{x}}_{n,k}^{\text{dry}}$ estimated by DDRM. After obtaining the estimates, we further refine the filter coefficients as

$$\mathbf{g}_k \leftarrow \mathbf{g}_k - \alpha \nabla_{\mathbf{g}_k} \left(\|\tilde{\mathbf{x}}_{n,k}^{\text{dry}} - (\tilde{\mathbf{I}} - \mathbf{G}_k) \mathbf{y}_{n,k}^{\text{wet}}\|_2^2 + \lambda \|\mathbf{g}_k\|_2^2 \right), \quad (10)$$

where α is a step size and λ is a regularization parameter. Note that \mathbf{G}_k is parameterized by \mathbf{g}_k and the derivative with respect to \mathbf{g}_k is tractable. With the refined $\mathbf{g}_k^{\text{refined}}$, the DDRM procedure is executed again to generate higher-quality results $\tilde{\mathbf{x}}_{n,k}^{\text{dry}}$. The proposed dereverberation algorithm is summarized in Fig. 1 and Algorithm 1.

4. EXPERIMENTS

4.1. Dataset

To examine the effectiveness of the proposed method, we conducted both quantitative and subjective evaluations on wet vocal signals. The pre-trained diffusion model was trained with only dry vocal signals from the NHSS dataset [26], which contains 100 English pop songs (20 unique songs) of ten different male and female singers. The total signal duration is 285.24 minutes.

As a test dataset, we prepared 3600 wet vocal signals, (ten hours in total) by adding artificial reverb to dry vocal signals from another dataset called the NUS-48E corpus [27]. This corpus contains 48 English pop songs (20 unique songs) of different male and female singers. The total signal duration is 169 minutes. Each song for both training and testing is sampled at 44.1 kHz and features monaural recording. As artificial reverb, we used all the presets for vocal in the FabFilter Pro-R plug-in¹. There are 19 kinds of vocal reverb presets in total. We prepared wet vocal signals by first making 48×19 wet vocal signals and dividing them into 10-second samples, and then randomly selecting 3600 signals from among them.

4.2. Experimental settings

The implementation of our method and the network architecture of the pre-trained diffusion model were mostly based on the code provided by the authors of the DDRM paper². We slightly modified some parts as follows. We converted each audio input into a complex-valued STFT representation using a window size of 1024, a hop size of 256, and a Hann window. Further, we cut the direct current component of the input signals and input them as 2-channel 512 \times 512 image data to follow the original input configurations. The first channel corresponds to the real value and the second to the imaginary value. We modified the original U-Net [28] architecture of the pre-trained model used on DDRM by adding a time-distributed fully connected layer [29] to the last layer of every residual block. For the training, we reduced the size of the diffusion model to have the fewer trainable parameters (31.3 M), and the training took less than three days using one NVIDIA A100 GPU.

For the inference, the parameters of WPE, DDRM, and the proposed refinement were set as follows. For WPE, we set $L = 150$ and $D = 8$, with the number of iterations set to one. For DDRM, we followed the same notations defined as in [5, Eq. (8)] and set $\eta = 0.7$, $\eta_b = 0.2$, and $\sigma_y = 1.0 \times 10^{-6}$, with the number of sampling steps set to 20. For the proposed refinement, we set $\alpha = 1.0 \times 10^{-6}$, $\lambda = 1.0$, and $N_{\text{refine}} = 10000$.

In addition, to explore the limitation of DDRM and our methods, we tested cases where a reverberation operator was able to be approximated from an oracle dry signal of test data $\mathbf{x}_{n,k}^{\text{test dry}}$ for a given wet observation of test data $\mathbf{y}_{n,k}^{\text{test wet}}$. We obtained this operator by minimizing $\|\mathbf{y}_{n,k}^{\text{test wet}} - \mathbf{H}_k \mathbf{x}_{n,k}^{\text{test dry}}\|_2^2$ with respect to \mathbf{H}_k over the reverberation operator.

4.3. Baselines

We evaluated the proposed method against three baselines.

Reverb conversion (RC): A state-of-the-art end-to-end DNN-based method for vocal dereverberation. We used the original code and the

¹<https://www.fabfilter.com/products/pro-r-reverb-plug-in>

²<https://github.com/bahjat-kawar/ddrm>

Table 1. Results of quantitative evaluation. ℓ_1 loss means ℓ_1 loss of magnitude spectrogram in STFT domain. Proposed and Proposed+ denote our methods without and with proposed refinement. DDRM w/ est-Oracle denotes case where DDRM was given an approximated oracle reverberation operator.

Methods	Manner	ℓ_1 loss ↓	FAD ↓
Wet (Unprocessed)	–	0.114	13.7
RC [1]	Supervised	0.117	13.9
ME [15]	Supervised	0.484	14.7
WPE [8]	Unsupervised	0.103	10.1
Proposed	Unsupervised	0.102	9.85
Proposed+	Unsupervised	0.100	9.69
DDRM w/ est-Oracle	–	0.079	4.44

pre-trained model³, which was trained with the pairs of 44.1 kHz wet and dry vocal signals. Note that the wet signals were reverbed with artificial reverb taken from the different professional reverb plug-ins from those of our test dataset [1]. We input pairs of wet and dry signals since this method needs them for dereverberation.

Music enhancement (ME): A supervised method to denoise and dereverb music signals based on diffusion models [15]. We used both the original code and the pre-trained model specified in the paper. Since ME was trained with pairs of 16 kHz reverberant noisy and clean music signals containing vocal signals, we evaluated this method at 16 kHz.

WPE: An unsupervised method for speech dereverberation [8]. We set $L = 200$, $D = 8$, and the number of iterations to three.

4.4. Quantitative evaluation

We evaluated the dereverberation performance of the proposed method without the refinement of an initial reverberation operator (Proposed) and with it (Proposed+) by computing two objective metrics. The first was the ℓ_1 loss of the amplitude spectrogram in the STFT domain between a dereverbed and dry signal. The other metric was the Fréchet audio distance (FAD) [30] between the set of dry and dereverbed signals. Since the VGGish [31], which is the pre-trained classification model of FAD, is originally trained with the 16 kHz audio samples, we downsampled all the signals to 16 kHz and computed FAD.

Table 1 lists the scores of each measurement. Both versions of the proposed method showed better scores than all the baselines on both metrics, and Proposed+ scored better than Proposed. These results demonstrate the effectiveness of our method and that utilizing the proposed refinement leads to a better dereverberation performance. The large gap between the scores of Proposed and DDRM w/ est-Oracle indicates that the estimation accuracy of reverberation operators significantly affects the dereverberation performance of our method. Thus, a better estimation of the initial reverberation operators and a better refinement of them will lead to a better dereverberation performance, which we leave to future work.

RC and ME did not perform well at all, which may according to that the distribution of their training dataset did not cover that of test dataset. Indeed, the training wet signals of ME and RC were created using only simulated natural reverb with some background noise [15] and different artificial reverb plug-ins from those of our test dataset [1], respectively. Another reason RC did not work well is

³The original code and the pre-trained model were shared by Junghyun Koo from the Artificial Intelligence Institute at Seoul National University. Mr. Koo also assisted with the discussion of the RC results of our experiment.

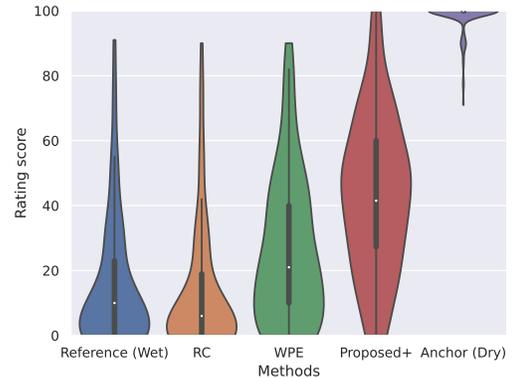


Fig. 2. Violin plots of listening test scores. White dots denote median of scores and top and bottom of vertical bold lines denote first and third quartiles, respectively.

that this model may not have been trained with enough pairs of wet and dry signals. Since RC was originally meant for not only dereverberation but also converting styles of reverb, it was also trained with the pairs of two different wet signals (see [1, Section 4.1]).

4.5. Listening test

We also conducted a listening test using the multiple stimuli with hidden reference and anchor (MUSHRA) method [32]. A total of 18 participants took part in the test on the webMUSHRA platform [33]. The participants were presented with ten kinds of signals (one for the practicing part and nine for the test part) selected randomly from the test dataset in Section 4.4 and trimmed to six seconds each. The results of the practicing part were removed from the aggregation of the results. Each web page contained a wet vocal signal as a reference and participants were asked to rate five different signals according to how much they felt the reverberant components were removed compared to the reference. The five signals were composed of the outputs of RC, WPE, Proposed+, a reference wet signal as a hidden reference, and a dry vocal signal as a hidden anchor. As with the results in Section 4.4, the outputs of ME and Proposed were not included here considering the burden on participants.

Figure 2 shows violin plots of the listening test results. Dry signals showed the highest score, which confirms that participants were able to judge which were wet and which were dry, as expected. Proposed+ showed the best score, outperforming both RC and WPE, which demonstrates the effectiveness of the proposed method, in terms of the perceptual metric.

5. CONCLUSION

In this work, we proposed an unsupervised method that can remove general reverb for music including artificial reverb without requiring diverse pairs of data for training. We extended DDRM to a practical use case with unknown reverberation operators. First, we initialize the reverberation operator with an estimation from WPE. Second, we adaptively refined the initial reverberation operator estimated by WPE to get a more accurate one with the help of a dry signal predicted by a diffusion model. The results of both objective and perceptual evaluations demonstrate that our method outperforms the current leading vocal dereverberation benchmarks.

6. REFERENCES

- [1] J. Koo, S. Paik, and K. Lee, “Reverb conversion of mixed vocal tracks using an end-to-end convolutional deep neural network,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2021, pp. 81–85.
- [2] V. Valimaki, J. D. Parker, L. Savioja, J. O. Smith, and J. S. Abel, “Fifty years of artificial reverberation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 20, no. 5, pp. 1421–1448, 2012.
- [3] T. Wilmering, D. Moffat, A. Milo, and M. B. Sandler, “A history of audio effects,” *Applied Sciences*, vol. 10, no. 3, 2020.
- [4] U. Zölzer, *DAFX: Digital Audio Effects: Second Edition*, WILEY, 03 2011.
- [5] B. Kawar, M. Elad, S. Ermon, and J. Song, “Denoising diffusion restoration models,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [6] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2020, vol. 33, pp. 6840–6851.
- [7] Y. Song, J. S.-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, “Score-based generative modeling through stochastic differential equations,” in *Proc. International Conference on Learning Representation (ICLR)*, 2021.
- [8] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, and B.-H. Juang, “Speech dereverberation based on variance-normalized delayed linear prediction,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, pp. 1717–1731, 2010.
- [9] N. Yasuraoka, T. Yoshioka, T. Nakatani, A. Nakamura, and H. G. Okuno, “Music dereverberation using harmonic structure source model and wiener filter,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2010, pp. 53–56.
- [10] N. Yasuraoka, H. Kameoka, T. Yoshioka, and H. G. Okuno, “I-divergence-based dereverberation method with auxiliary function approach,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2011, pp. 369–372.
- [11] K. Mahkonen, A. Eronen, T. Virtanen, E. Helander, V. Popa, and I. D. D. Curcio, “Music dereverberation by spectral linear prediction in live recordings,” in *Proceedings of International Conference on Digital Audio Effects*, 2013.
- [12] A. Maezawa, K. Itoyama, K. Yoshii, and H. G. Okuno, “Non-parametric bayesian dereverberation of power spectrograms based on infinite-order autoregressive processes,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 22, pp. 1918–1930, 2014.
- [13] T. Okamoto, Y. Iwaya, and Y. Suzuki, “Wide-band dereverberation method based on multichannel linear prediction using prewhitening filter,” *Applied Acoustics*, vol. 73, no. 1, pp. 50–55, 2012.
- [14] A. Tsilfidis and J. Mourjopoulos, “Blind single-channel dereverberation for music post-processing,” *J. Audio Eng. Soc.*, 2011.
- [15] N. Kandpal, O. Nieto, and Z. Jin, “Music enhancement via image translation and vocoding,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2022, pp. 3124–3128.
- [16] A. Ribes and F. Schmitt, “Linear inverse problems in imaging,” *IEEE Signal Process. Mag.*, vol. 25, no. 4, pp. 84–99, 2008.
- [17] A. Adler, V. Emiya, M. G. Jafari, M. Elad, R. Gribonval, and M. D. Plumbley, “Audio inpainting,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 20, no. 3, pp. 922–932, 2011.
- [18] S. Arberet, P. Vandergheynst, R. E. Carrillo, J.-P. Thiran, and Y. Wiaux, “Sparse reverberant audio source separation via reweighted analysis,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 21, no. 7, pp. 1391–1402, 2013.
- [19] D. Ulyanov, A. Vedaldi, and V. Lempitsky, “Deep image prior,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 9446–9454.
- [20] V. Shah and C. Hegde, “Solving linear inverse problems using gan priors: An algorithm with provable guarantees,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2018, pp. 4609–4613.
- [21] Z. Kadhodaie and E. Simoncelli, “Stochastic solutions for linear inverse problems using the prior implicit in a denoiser,” *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, pp. 13242–13254, 2021.
- [22] B. Kawar, G. Vaksman, and M. Elad, “SNIPS: Solving noisy inverse problems stochastically,” *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, pp. 21757–21769, 2021.
- [23] H. Chung, B. Sim, D. Ryu, and J. C. Ye, “Improving diffusion models for inverse problems using manifold constraints,” *arXiv preprint arXiv:2206.00941*, 2022.
- [24] H. Chung, B. Sim, and J. C. Ye, “Come-closer-diffuse-faster: Accelerating conditional diffusion models for inverse problems through stochastic contraction,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 12413–12422.
- [25] C.-H. Lai, Y. Takida, N. Murata, T. Uesaka, Y. Mitsufuji, and S. Ermon, “Regularizing score-based models with score fokker-planck equations,” *arXiv preprint arXiv:2210.04296*, 2022.
- [26] B. Sharma, X. Gao, K. Vijayan, X. Tian, and H. Li, “NHSS: A speech and singing parallel database,” *arXiv preprint arXiv:2012.00337*, 2020.
- [27] Z. Duan, H. Fang, B. Li, and Y. Sim, K. C. and Wang, “The NUS sung and spoken lyrics corpus: A quantitative comparison of singing and speech,” in *Proc. Asia-Pacific Signal Inf. Process. Assoc. Annu. Summit Conf. (APSIPA ASC)*, 2013, pp. 1–9.
- [28] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *Proceedings of the International Conference on Medical image computing and computer-assisted intervention*, 2015, pp. 234–241.
- [29] W. Choi, M. Kim, J. Chung, D. Lee, and S. Jung, “Investigating U-Nets with various intermediate blocks for spectrogram-based singing voice separation,” in *Proc. Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2020.
- [30] K. Kilgour, M. Zuluaga, D. Roblek, and M. Sharif, “Fréchet Audio Distance: A metric for evaluating music enhancement algorithms,” *arXiv preprint arXiv:1812.08466*, 2018.
- [31] S. Hershey, S. Chaudhuri, D. P. W. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold, M. Slaney, R. J. Weiss, and K. Wilson, “Cnn architectures for large-scale audio classification,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2017, pp. 131–135.
- [32] I. T. Union, “Recommendation itu-r bs.1534-1: Method for the subjective assessment of intermediate quality level of coding systems,” 2003.
- [33] M. Schoeffler and et al., “webMUSHRA — A comprehensive framework for web-based listening tests,” *A Comprehensive Framework for Web-based Listening Tests. Journal of Open Research Software.*, vol. 6, no. 1, pp. 8, 2018.
- [34] P. Dhariwal and A. Nichol, “Diffusion models beat GANs on image synthesis,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2021, vol. 34, pp. 8780–8794.
- [35] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *Proc. International Conference on Learning Representation (ICLR)*, 2019.
- [36] P. Micikevicius, S. Narang, J. Alben, G. Diamos, E. Elsen, D. Garcia, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh, and H. Wu, “Mixed precision training,” in *Proc. International Conference on Learning Representation (ICLR)*, 2018.
- [37] Y. Song and S. Ermon, “Improved techniques for training score-based generative models,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2020, vol. 33, pp. 12438–12448.
- [38] A. Nichol and P. Dhariwal, “Improved denoising diffusion probabilistic models,” *arXiv preprint arXiv:2102.09672*, 2021.

7. APPENDIX

In this section, we additionally demonstrate the detailed settings of the training configurations of the pre-trained diffusion model. The network architecture of the pre-trained diffusion model is mostly based on the code provided by the authors of the DDRM paper ⁴. Especially, we explain our modifications about the input representation and the network architecture in Section 4.2. The pre-trained diffusion models in the DDRM code are from this *guided-diffusion* GitHub repository ⁵ [34]. The hyperparameters for the training of the diffusion model are in Table 2. We also incorporate an adaptive group normalization [34] into each residual block. We train the model using AdamW [35] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ in 16-bit precision [36]. We use an exponential moving average over model parameters with a rate of 0.9999 [37].

Table 2. Hyperparameters for training diffusion model. We followed same notations defined in [34, Table 11]

Diffusion steps	4000
Noise schedule	cosine [38]
Model size	31.3 M
Channels	64
Depth	2
Channels multiple	1, 1, 2, 2, 4, 4
Heads	2
Attention resolution	32, 16
BigGAN up/downsample	✓
Dropout	0.0
Batch size	6
Iterations	370K
Learning rate	1.0×10^{-4}

⁴<https://github.com/bahjat-kawar/ddrm>

⁵<https://github.com/openai/guided-diffusion>