



HAL
open science

Diffusion-based spectral super-resolution of third octave acoustic sensor data: is privacy at risk ?

Modan Tailleur, Chaymae Benaatia, Mathieu Lagrange, Pierre Aumond, Vincent Tourre

► **To cite this version:**

Modan Tailleur, Chaymae Benaatia, Mathieu Lagrange, Pierre Aumond, Vincent Tourre. Diffusion-based spectral super-resolution of third octave acoustic sensor data: is privacy at risk?. 33rd European Signal Processing Conference (EUSIPCO 2025), European Association for Signal Processing (EURASIP), Sep 2025, Palerme, Italy. pp.306. <hal-05096000v2>

HAL Id: hal-05096000

<https://hal.science/hal-05096000v2>

Submitted on 21 Jul 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Copyright - All rights reserved

Diffusion-based Spectral Super-Resolution of Third Octave Acoustic Sensor Data: Is Privacy at Risk ?

1st Modan Tailleux*, 2nd Chaymae Benaatia*, 3rd Mathieu Lagrange*, 4th Pierre Aumond†, 5th Vincent Tourre‡

*Nantes Université, École Centrale Nantes, CNRS, LS2N, UMR 6004, F-44000 Nantes, France

†Université Gustave Eiffel, CEREMA, UMRAE, F-44344 Bouguenais, France

‡Nantes Université, ENSA Nantes, École Centrale Nantes, CNRS, AAU-CRENAU, UMR 1563, F-44000 Nantes, France

Abstract—Third octave spectral recording of acoustic sensor data is an effective way of measuring the environment. While there is strong evidence that slow (1s frame, 1 Hz rate) and fast (125ms frame, 8Hz rate) versions lead by-design to unintelligible speech if reconstructed, the advent of high quality reconstruction methods based on diffusion may pose a threat, as those approaches can embed a significant amount of *a priori* knowledge when learned over extensive speech datasets.

This paper aims to assess this risk at three levels of attacks with a growing level of *a priori* knowledge considered at the learning of the diffusion model, a) none, b) multi-speaker data excluding the target speaker and c) target speaker. Without any prior regarding the speech profile of the speaker (levels a and b), the word-error-rate both for humans and automatic recognition remains higher than 89%.

Index Terms—content privacy, generative audio, acoustic sensor networks, audio encoding

I. INTRODUCTION

In recent years, the use of acoustic sensors for audio data collection has extended across diverse applications, encompassing domains such as smart homes [1], [2] and urban sensor networks [3], [4]. Ensuring the privacy of speech information is a key aspect of the deployment of such sensors, be it deployed on public or private places. A promising approach that emerged from previous studies involves the encoding of audio as fast third-octave spectrograms (FTOS), which are third-octave spectro-temporal data computed with 125ms windows and 20-29 frequency bands [5]. By considering a low sampling rate, this method proved effective in preserving content privacy as phoneme average duration in spoken English is typically below 100ms [6]. With longer windows, the co-articulation of phonemes is lost, leading to an almost complete loss of intelligibility. This loss ensures “content privacy”, simply termed in this paper “privacy”. Speaker privacy by means of its identity is another important issue that is not considered in this paper.

Considering an at the time state-of-the-art reconstruction approach combining a) the Moore-Penrose pseudo-inverse (PINV) for frequency retrieval and b) the Griffin-Lim algorithm [7] for phase retrieval, Gontier et al. empirically

demonstrated that the recovered speech was unintelligible [8]. Based purely on signal processing techniques, this reconstruction method does not consider any *a priori* knowledge on spectro-temporal properties of spoken English. However, to our knowledge, no attempts have been made to recover speech information from FTOS data using deep learning methods. With the recent advancements in generative audio models [9], which leverage *a priori* knowledge from large amounts of speech data, we believe that there is a need to re-evaluate the aforementioned claim that is: FTOS encoding is *by-design* preserving content privacy.

Particularly, the emergence of diffusion models [10] may pose a threat. These models are easier to train than Generative Adversarial Networks (GANs) and can thus be applied to a broader range of fields. Indeed, diffusion models have demonstrated super-resolution abilities [11], [12], and have notably shown good performances in enhancing the quality of speech [13], [14].

In order to evaluate potential privacy threats induced by training such algorithms, we categorize out attack models into three distinct Attack Levels (AL) on FTOS data, based on training set selection:

- **AL0** denotes an unintentional attack, occurring when a model is trained on general urban data or without any specialized training. The aim of the attacker in this case is not to recover specifically speech information but rather to reconstruct general audio for analysis purposes.
- **AL1** denotes an attack resulting from training a model on general speech data. In this scenario, the attacker aims to recover speech information from an unknown speaker.
- **AL2** denotes an attack executed by training a model specifically on the voice of a target individual. In this scenario, it is assumed that the attacker has access to a sufficient amounts of clear speech data from his target.

To evaluate the impact of large-scale attacks on sensor data, we anticipate that attackers will utilize Automatic Speech Recognition (ASR) systems on FTOS data. Given the uncertainty surrounding the effectiveness of ASRs in processing severely downsampled acoustic signals such as FTOS, we perform a subjective assessment through Human Speech Recognition (HSR) evaluation. This subjective evaluation provides a

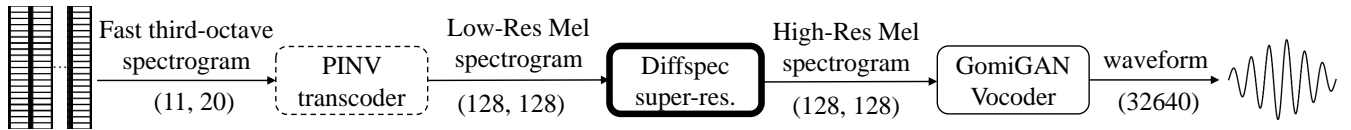


Fig. 1. Pipeline for audio super-resolution. Module outlined with dashes does not require any training, module with a plain outline is only pre-trained, and module outlined in bold is specifically trained for the task.

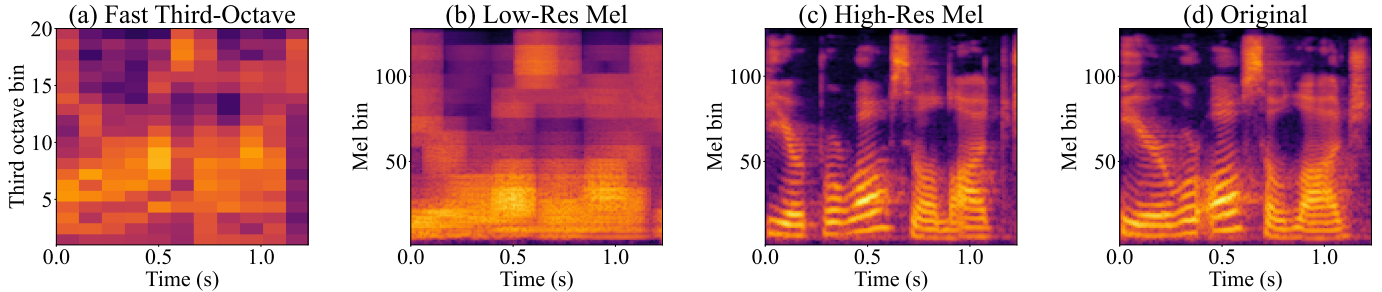


Fig. 2. Log spectrograms of the LJ002-0068.wav audio file from the LJSpeech evaluation set. The DiffSpec model, trained on the LJSpeech training set, reconstructs a High-Res spectrogram (c) that closely matches the original (d).

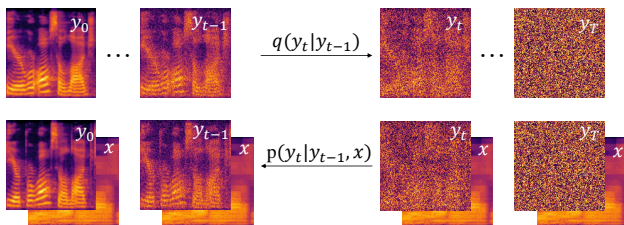


Fig. 3. Forward process (top): Gaussian noise progressively added to the target Mel spectrogram. Backward process (bottom): generative denoising conditioned by Low-Res Mel spectrogram x on second channel.

	FTOS	GomiGAN Mel
sample rate	32kHz	24kHz
window size	4096 (128ms)	1024 (43ms)
hop size	4000 (125ms)	320 (11ms)
window	Tukey	Hann
frequency bins	20	128
min frequency	125Hz	23Hz
max frequency	10kHz	12kHz

TABLE I
DIFFERENCES BETWEEN FAST THIRD-OCTAVE SPECTROGRAMS (FTOS)
AND GOMIGAN MEL INPUTS.

reference for interpreting ASRs performance under those kind of adverse conditions.

In Section II, we present a super resolution technique based on diffusion [12] to recover speech from FTOS encoded speech. In Section III, we describe the experimental protocol used to assess the privacy threat posed by this method. Sections IV and V, detail our main findings based on perceptual and computational assessment of intelligibility under scenarios simulating the 3 levels of attack¹.

II. METHOD

Our proposed approach, which we refer to as the **DiffSpec** method, transform FTOS into Mel spectrograms using a super-resolution algorithm, and use a vocoder for phase reconstruction. As audio is recovered, off-the-shelf ASRs or HSR can be used to recover speech information. We believe this approach is effective because it minimizes the need for extensive training. By relying on a pre-trained vocoder for the

vocoding stage, the model only needs to focus on spectrogram reconstruction, thereby simplifying the overall process.

As shown in Figure 1, we choose to use a DiffSpec pipeline considering Saharia et al. [12] super-resolution algorithm applied on spectrograms, a conditioning method notably utilized in the NU-Wave 2 algorithm [15]. For the vocoder, we select the pre-trained GomiGAN model [16]. Consequently, the pipeline must align with the input dimensions of the GomiGAN model, as shown in Table I. Going from FTOS to GomiGAN mels requires a 74x upscaling factor. After the vocoding stage, audio segments of 1.36s are computed and then concatenated with hops of 1.23s and cross-fades to match the initial audio length. An example of the spectrograms generated in the different DiffSpec stages is shown in Figure 2.

First, we create an initial low-resolution approximation of a GomiGAN Mel spectrogram from the FTOS. We use a Moore-Penrose pseudo-inverse followed by Mel filtering to align with the target Mel frequency bins, and linear interpolation to match the number of time frames. This spectrogram will be referred to as **Low-Res Mel spectrogram** in the following sections. Detailed methodology for obtaining this pseudo-inverted Mel spectrogram can be found in Tailleir et al. [17].

¹Code operating solely on public data and audio examples are available at: <https://modantailleir.github.io/paperThirdOctavePrivacy/>

The Low-Res Mel spectrogram is then refined using a DiffSpec based on the model proposed by Saharia et al. [12]. This refined spectrogram will be referred to as **High-Res Mel spectrogram**. The U-net model used for the diffusion includes two ResNet layers per block, organized into six blocks with 64, 64, 128, 128, 256, and 256 output channels, respectively, and contains a total of 28 million parameters. It features one attention block in both the downstream and upstream stages. The forward and backward diffusion processes, using second-channel conditioning, are presented in Figure 3.

We then apply the GomiGAN [16] vocoder on the output of the DiffSpec model. GomiGAN is a general-purpose vocoder designed to convert any Mel spectrogram with characteristics shown in Table I to waveform audio. It is trained on a diverse range of audio datasets, including speech signals, music stems, animal sound recordings, and foley sounds. The GomiGAN model is based on BigVGAN [18], enhanced with Feature-wise Linear Modulation (FiLM) [19]. While the Griffin-Lim algorithm is also a potential vocoder alternative, GomiGAN offers several advantages. It runs approximately 15 times faster than Griffin-Lim with 32 iterations due to its GPU compatibility, and informal listening by the authors indicated superior audio quality. Objective and subjective evaluation have been considered for both vocoders. Only the DiffSpecs using GomiGAN as vocoder are reported in Section V, as no significant differences between the two vocoders are found.

III. EXPERIMENTAL PROTOCOL

Data

For each Attack Level (AL) defined in Section I, we select a specific audio dataset to compute FTOS data to train our model. The chosen datasets are:

- **AL0** dataset: TAU Urban Acoustic Scenes 2020 Mobile dataset [20]. This dataset contains 10-second audio clips from 10 different acoustic scenes, including indoor public spaces, public transports, streets, and parks. While it includes some distant voice samples, it primarily focuses on ambient urban sounds, and totals 64 hours of audio.
- **AL1** dataset: Librispeech [21]. It consist of read audio-books with more than 2,000 speakers. Specifically we use the "train-clean-100" subset, which includes 100 hours of audio data.
- **AL2** dataset: LJSpeech ². It comprises 13,000 audio clips from a single speaker reading seven non-fiction books, totalling 24h of audio. Like Librispeech, readings are available through the LibriVox project. The training dataset we considered comprises 12,900 audio clips, as 100 are kept for evaluation.

For evaluation, we randomly select 100 audio samples from the LJSpeech dataset. The relatively small size of this subset is due to the high computational cost of diffusion model inference, which can take up to 50 seconds per 10-second audio sample on a GeForce RTX 2080 Ti GPU.

Baseline

²LJSpeech dataset available at: <https://keithito.com/LJ-Speech-Dataset/>

We compare our model against a simple pseudo-inverse approach using a PINV transcoder, as described in section II. Compared to the proposed approach, this baseline simply bypasses the DiffSpec step and applies the Vocoder directly to the Low-Res Mel spectrogram.

Learning procedure

The DiffSpec model is trained with a learning rate of 10^{-4} , a batch size of 200, for 40,000 iterations.

Metric

To assess the privacy threat potentially induced by the generated audio samples, we measure the Word Error Rate (WER). The WER is a measure of the discrepancy between the reference transcriptions and those produced by HSR or ASR systems on the reconstructed speech. It is calculated as the sum of the number of substitutions, deletions, and insertions required to convert the inferred text into the reference text, divided by the total number of words in the reference text.

IV. SUBJECTIVE EVALUATION

A subjective evaluation is performed on audio data reconstructed from FTOS, using WER on transcriptions from fluent english speakers who have reported normal hearing. Before the final analysis, the first author manually performs obvious grammar and typos corrections on all participants transcriptions.

8 audio samples are selected from the 100 audio samples of our LJSpeech evaluation subset, and are transformed into FTOS. These samples are chosen to be at least 6-s long and to contain content understandable without extensive cultural knowledge, avoiding names and slangs.

From informal listening done by the authors, some of the settings obviously lead to either full unintelligibility (AL0) or full intelligibility (original mel processed through GomiGAN). To keep the final perceptual test tractable and avoid cluttering the evaluation with settings that show highly contrasting WERs, we decide to evaluate those highly contrasted settings on a initial test conducted with only 3 participants.

Our 20 other participants transcribe audios generated from the remaining two settings, which are the DiffSpec models trained on Librispeech and LJSpeech (AL1 and AL2). Each participant transcribes a total of 16 audio samples: 4 samples from each of the two systems (AL1 and AL2) and 8 from the original audios. Participants whose transcriptions of the original audios lead to a WER exceeding 10% are excluded from the analysis. As a result, 3 participants are removed, leaving data of 17 participants for analysis.

As shown in Table II, our reference human speech recognition (HSR) leads to a minimum of 92% WER on AL0, 90% of WER on AL1 and 64% of WER on AL2.

V. OBJECTIVE EVALUATION

An objective evaluation is then performed using the WER computed for different state-of-the-art automatic speech recognition (ASR) models: Wav2Vec2 [22], the "large-v3" Whisper model [23], Fairseq S2T [24], as well as the CRDNN model

Attack Level	Training Set	Method	HSR	FairseqS2T	W2V2	CRDNN	Whisper
-	-	Original (mel)	01 (± 01)	10 (± 03)	10 (± 03)	08 (± 03)	02 (± 01)
-	-	White Noise	-	97 (± 01)	100 (± 00)	99 (± 01)	95 (± 01)
AL0	-	PINV transc.	98 (± 01)	97 (± 01)	95 (± 01)	98 (± 01)	82 (± 04)
	TAU	Diffspec	92 (± 03)	94 (± 02)	93 (± 02)	94 (± 01)	92 (± 03)
AL1	Librispeech	Diffspec	90 (± 02)	91 (± 02)	89 (± 02)	91 (± 02)	85 (± 03)
AL2	LJSpeech	Diffspec	64 (± 04)	53 (± 05)	52 (± 04)	46 (± 04)	35 (± 04)

TABLE II

WORD ERROR RATE (IN %) ON LJSPEECH FOR THE DIFFERENT COMBINATIONS OF METHODS AND TRAINING SETS. HSR REPRESENTS THE HUMAN SPEECH RECOGNITION EVALUATED WITH PERCEPTUAL EXPERIMENT. THE "ORIGINAL (MEL)" METHOD DESIGNATES THE ORIGINAL AUDIO TRANSFORMED INTO A GOMIGAN MEL SPECTROGRAM AND RAN THROUGH GOMIGAN VOCODER. THE CONFIDENCE INTERVAL SHOWN IS A 95% CONFIDENCE INTERVAL.

from the Speechbrain library [25] called "asr-crdnn-rnnlm-librispeech". Table II presents the results of this evaluation.

The results indicate that all ASR systems yield WERs comparable to HSR, with the exception of Whisper. As they even outperform human evaluations at the AL2 attack level, this suggests that ASR systems are generally robust and effective for processing reconstructed audio. The Whisper ASR system exhibits a rather peculiar behavior, with overall low WERs across all attack levels. This is especially worrying in the case of the PINV transcoder, showing a surprisingly low 82% WER, where human listening demonstrates almost complete unintelligibility (audio examples are provided on the companion page). We suspect that this behavior may be due to some sort of overfitting of this specific ASR to the LJSpeech dataset. Due to a lack of full understanding, the performance of Whisper is not discussed further.

For an unintentional attack resulting from training on a non-speech dataset or using a non-learned algorithms (attack level AL0), the results from the diverse ASR systems show that speech is nearly unintelligible. In this scenario, the WER is only a few percentage points lower than the one obtained with white noise, which is consistent with the results of Gontier et al. [8].

When targeting specifically speech information (attack level AL1), the results are slightly more concerning. W2V2 notably achieves an 89% WER. Although this might not seem alarming, understanding 10% of words could pose significant privacy risks in certain contexts, particularly when deploying systems in sensitive or private environments. Interestingly, the recovered words are not limited to common function words (e.g., "and", "or") but sometimes include more semantically meaningful content. Moreover, the model appears particularly effective at recovering full words when the speaker articulates very slowly.

Targeting not only speech information but specifically the voice of the LJSpeech speaker (attack level AL2), the WER on W2V2 reaches up to 46%. While this level of WER indicates that a significant portion of the speech remains unintelligible, it is important to recognize that comprehending more than 50% of the words in a conversation might allow to grasp the overall meaning of the sentences. However, this scenario remains extreme, as it only occurs when a model is specifically trained on the target speaker's voice.

We also conducted experiments on a Librispeech evaluation

set for attack levels AL0 and AL1, as AL2 could not be performed on this dataset. Due to space constraints, and because the trends were consistent with those observed on LJSpeech evaluation set, the full table is provided only on the companion page.

VI. CONCLUSION

Using the Diffspec method across the different Attack Levels (AL) we have established, our model demonstrates a minimum Word Error Rate (WER) of 93% for AL0, 89% for AL1, and 46% for AL2 for several ASR systems on our LJSpeech test set.

AL0 represents a common use case for urban sound monitoring, as reconstructing waveforms from FTOS enables audio playback of the recorded data. The fact that the Attack Level performs only slightly better than White Noise suggests that using Diffspec, when not trained on speech, is generally safe for urban sound monitoring from a privacy perspective. Results on AL1 and AL2 also indicate a very low risk of extracting intelligible speech information when Diffspec is trained on speech data. However, even with the seemingly high WERs for AL1 and AL2, the risk associated with these attacks is highly context-dependent. Future work could therefore consider perceptual assessment of acceptable risk thresholds for WER in specific application settings from public spaces, to offices or bedrooms.

We have shown in this paper that with the rise of generative models, fast third-octave spectrograms are no longer inherently privacy-aware as initially suggested by Gontier et al. [8]. Future research could involve training and testing models on multiple individual voices to further validate and strengthen these findings.

We believe that the proposed Diffspec approach nicely balances audio quality and training requirements in terms of data and power. By enabling reconstruction in the waveform domain, it also enhances explainability, allowing for direct comparison between the results of ASR and HSR. While Diffspec performs well in our experiments, other generative architectures may yield different outcomes worth exploring. Furthermore, more complex approaches could be considered. For example, one could train an end-to-end automatic speech recognition (ASR) algorithm that uses fast third-octave spectrograms (FTOS) as input instead of Mel spectrograms [26]–[29], though this approach would require intensive training

to reach the performance of off-the-shelf ASRs. One could also consider training vocoders to convert third-octave spectrograms directly to waveforms [30]–[33]. This latter approach could prove to be useful, but preliminary attempts by the authors demonstrated that the size of training dataset and computational power needed for training is notably larger than the ones required by the DiffSpec method.

Finally, preliminary experiments applying DiffSpec to non-speech audio have shown promising results for general super-resolution of environmental sound scenes under ALO. Future work should further evaluate its performance in this domain, compared to the PINV model.

REFERENCES

- [1] Sacha Krstulović, “Audio Event Recognition in the Smart Home,” in *Computational Analysis of Sound Scenes and Events*. 2018.
- [2] Michel Vacher, Dan Istrate, François Portet, Thierry Joubert, Thierry Chevalier, Serge Smidtas, Brigitte Meillon, Benjamin Lecouteux, Mohamed Sehili, and Pedro Chahua, “The sweet-home project: Audio technology in smart homes to improve well-being and reliance,” in *EMBC*, 2011.
- [3] Arnaud Can, Judicaël Picaut, Jeremy Ardouin, Pierre Crepeaux, Erwan Bocher, David Ecotiere, and Mathieu Lagrange, “CENSE Project: general overview,” in *EuroNoise*, 2021.
- [4] Mark Cartwright, Ana Elisa Mendez Mendez, Jason Cramer, Vincent Lostonlen, Graham Dove, Ho-Hsiang Wu, Justin Salamon, Oded Nov, and Juan Bello, “SONYC urban sound tagging (SONYC-UST): A multilabel dataset from an urban acoustic sensor network,” in *DCASE*, 2019.
- [5] Jérémy Ardouin, Ludovic Charpentier, Mathieu Lagrange, Félix Gontier, Nicolas Fortin, David Ecotière, Judicaël Picaut, and Christophe Mietlicky, “An innovative low cost sensor for urban sound monitoring,” in *InterNoise*, 2018.
- [6] Simon W McKnight, Aidan OT Hogg, and Patrick A Naylor, “Analysis of phonetic dependence of segmentation errors in speaker diarization,” in *EUSIPCO*, 2021.
- [7] Daniel Griffin and Jae Lim, “Signal estimation from modified short-time Fourier transform,” *IEEE TASLP*, 1984.
- [8] Félix Gontier, Mathieu Lagrange, Pierre Aumond, Arnaud Can, and Catherine Lavandier, “An efficient audio coding scheme for quantitative and qualitative large scale acoustic monitoring using the sensor grid approach,” *MDPI Sensors*, 2017.
- [9] Staphord Bengesi, Hoda El-Sayed, Md Kamruzzaman Sarker, Yao Houkpati, John Irungu, and Timothy Oladunni, “Advancements in generative ai: A comprehensive review of gans, gpt, autoencoders, diffusion model, and transformers,” *IEEE Access*, 2024.
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel, “Denoising diffusion probabilistic models,” *NeurIPS*, 2020.
- [11] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song, “Denoising diffusion restoration models,” *NeurIPS*, 2022.
- [12] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J. Fleet, and Mohammad Norouzi, “Image super-resolution via iterative refinement,” *IEEE TPAMI*, 2022.
- [13] Joan Serrà, Santiago Pascual, Jordi Pons, R. Oguz Araz, and Davide Scaini, “Universal speech enhancement with score-based diffusion,” 2022.
- [14] Ryosuke Sawata, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Takashi Shibuya, Shusuke Takahashi, and Yuki Mitsufuji, “Diffiner: A Versatile Diffusion-based Generative Refiner for Speech Enhancement,” in *Interspeech*, Aug. 2023.
- [15] Seungu Han and Junhyeok Lee, “Nu-wave 2: A general neural audio upsampling model for various sampling rates,” in *Interspeech*, 2022.
- [16] Minsung Kang, Sangshin Oh, Hyeonggi Moon, Kyungyun Lee, and Ben Sangbae Chon, “FALL-E: A Foley Sound Synthesis Model and Strategies,” in *DCASE*, 2023.
- [17] Modan Tailleir, Mathieu Lagrange, Pierre Aumond, and Vincent Tourre, “Spectral transcoder: using pretrained urban sound classifiers on under-sampled spectral representations,” in *DCASE*, 2023.
- [18] Sang-gil Lee, Wei Ping, Boris Ginsburg, Bryan Catanzaro, and Sungroh Yoon, “BigVGAN: A Universal Neural Vocoder with Large-Scale Training,” in *ICLR*, 2022.
- [19] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville, “Film: Visual reasoning with a general conditioning layer,” in *AAAI*, 2018.
- [20] Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen, “A multi-device dataset for urban acoustic scene classification,” in *DCASE*, 2018.
- [21] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *ICASSP*, 2015.
- [22] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *NeurIPS*, 2020.
- [23] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever, “Robust speech recognition via large-scale weak supervision,” in *ICML*, 2023.
- [24] Changhan Wang, Yun Tang, Xutai Ma, Anne Wu, Dmytro Okhonko, and Juan Pino, “Fairseq S2T: Fast Speech-to-Text Modeling with Fairseq,” in *AAACL-IJCNLP*, 2020.
- [25] Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, Loren Lugosch, Cem Subakan, Nauman Dawalatabad, Abdelwahab Heba, Jianyuan Zhong, Ju-Chieh Chou, Sung-Lin Yeh, Szu-Wei Fu, Chien-Feng Liao, Elena Rastorgueva, François Grondin, William Aris, Hwidong Na, Yan Gao, Renato De Mori, and Yoshua Bengio, “Speechbrain: A general-purpose speech toolkit,” 2021.
- [26] Sungwon Kim, Heeseung Kim, and Sungroh Yoon, “Guided-its 2: A diffusion model for high-quality adaptive text-to-speech with untranscribed data,” 2022.
- [27] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu, “Diffsound: Discrete Diffusion Model for Text-to-Sound Generation,” *IEEE TASLP*, 2023.
- [28] Jaesung Tae, Hyeongju Kim, and Taesu Kim, “EdiTTs: Score-based Editing for Controllable Text-to-Speech,” in *Interspeech*, 2021.
- [29] Rongjie Huang, Zhou Zhao, Huadai Liu, Jinglin Liu, Chenye Cui, and Yi Ren, “ProDiff: Progressive Fast Diffusion Model for High-Quality Text-to-Speech,” in *ACM MM*, Oct. 2022.
- [30] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro, “DiffWave: A Versatile Diffusion Model for Audio Synthesis,” in *ICLR*, 2021.
- [31] Tan Dat Nguyen, Ji-Hoon Kim, Youngjoon Jang, Jaehun Kim, and Joon Son Chung, “Fregrad: Lightweight and Fast Frequency-Aware Diffusion Vocoder,” in *ICASSP*, 2024.
- [32] Sang-gil Lee, Heeseung Kim, Chaehun Shin, Xu Tan, Chang Liu, Qi Meng, Tao Qin, Wei Chen, Sungroh Yoon, and Tie-Yan Liu, “PriorGrad: Improving Conditional Denoising Diffusion Models with Data-Dependent Adaptive Prior,” in *ICLR*, 2022.
- [33] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, Najim Dehak, and William Chan, “WaveGrad 2: Iterative Refinement for Text-to-Speech Synthesis,” in *Interspeech*, 2021.