

GENCHO: ROOM IMPULSE RESPONSE GENERATION FROM REVERBERANT SPEECH AND TEXT VIA DIFFUSION TRANSFORMERS

Jackie Lin^{1,2}, Jiaqi Su², Nishit Anand^{2,3}, Zeyu Jin², Minje Kim¹, and Paris Smaragdis⁴

¹University of Illinois Urbana-Champaign, ²Adobe Research, ³University of Maryland, ⁴MIT

ABSTRACT

Blind room impulse response (RIR) estimation is a core task for capturing and transferring acoustic properties; yet existing methods often suffer from limited modeling capability and degraded performance under unseen conditions. Moreover, emerging generative audio applications call for more flexible impulse response generation methods. We propose Gencho, a diffusion-transformer-based model that predicts complex spectrogram RIRs from reverberant speech. A structure-aware encoder leverages isolation between early and late reflections to encode the input audio into a robust representation for conditioning, while the diffusion decoder generates diverse and perceptually realistic impulse responses from it. Gencho integrates modularly with standard speech processing pipelines for acoustic matching. Results show richer generated RIRs than non-generative baselines while maintaining strong performance in standard RIR metrics. We further demonstrate its application to text-conditioned RIR generation, highlighting Gencho’s versatility for controllable acoustic simulation and generative audio tasks.

Index Terms— blind room impulse response estimation, acoustic matching, generative AI, diffusion transformer

1. INTRODUCTION

Room impulse responses (RIRs) are filters that capture the core acoustic properties of an environment, including reverberation and coloration, through a compact, parametric representation. They provide essential auditory context that shapes how sound is perceived in a given space, contributing to the realism and immersiveness of audio content. RIR estimation provides a natural basis for acoustic matching, which aims to transfer the acoustics of a reference space to new audio so that it blends seamlessly with the original scene.

Recently, the proliferation of publicly available audio crafting tools, such as speech enhancement and text-to-speech (TTS) synthesis, has further increased the need for more accurate, flexible acoustic matching. This capability is critical for tasks such as automated dialogue replacement (ADR), dubbing, and voiceovers, to ensure perceptual consistency of newly recorded or synthetic speech with the original context. Likewise, TTS integrates with acoustic matching to render voices consistently in a chosen environment. Processed speech from enhancers can sound overly dry, and benefits from restoring natural room acoustics. In many scenarios, explicitly estimating RIRs is often preferred over end-to-end acoustic matching, as it produces reusable filters that can be stored, edited, shared, and applied across tasks without altering the underlying audio.

Moreover, with the emerging volume of generative tools and synthetic content, the scope of IR estimation is expanding beyond acoustic matching to a real audio reference. Applications such as

immersive storytelling, AR/VR, and text-to-audio generation require *soft* acoustic matching: the ability to generate diverse, semantically appropriate RIRs from weak or indirect cues—images, videos, or natural-language descriptions—in order to create coherent and immersive virtual acoustic environments. Here, generating acoustics separately from content allows users to modify acoustics without affecting speaker identity or other semantic properties and vice versa.

Despite much work, estimating RIRs from audio recordings in a blind setting—with no prior knowledge of the room, recording setup, or source signal—remains a fundamental challenge in audio processing. Lightweight parametric models are grounded in strong inductive biases, limiting their ability to capture the complexity of real RIRs. For instance, blind parameter estimation from reverberant speech [1, 2, 3] predicts octave-band reverberation times (T60) which can be used to drive DSP-based reverberators; however, the coarseness of these parameters limits the diversity of the synthesized RIRs. A more recent deep learning approach, Filtered Noise Shaping (FiNS) [4], constructs RIRs by upsampling a learned vector in time domain to produce an early reflection waveform and noise envelopes for the late tail. However, the output space remains constrained by the model formulation, often producing similar-sounding or unreasonable RIRs on unseen reverberant speech. In general, non-generative systems fail to capture the multi-modal nature of the blind inverse problem, and are prone to degeneracy when faced with unfamiliar inputs. As a response, GAN-based [5] and language modeling-based [6, 7] generative models have been introduced to mitigate these issues, yet fully closing the perceptual gap between synthesized and real acoustics remains an open challenge.

We introduce Gencho¹ (**Generative Echo**), a blind room impulse response estimator using diffusion transformers. We address the limitations of non-generative acoustic matching methods with our method’s strong generalization capability and ability to produce diverse, in-distribution RIRs. Moreover, Gencho is designed to work with modern audio technologies; by leveraging speech enhancement and source separation to extract dry and early reflected speech signals, our pipeline focuses specifically on RIR estimation while integrating seamlessly into end-to-end workflows. Lastly, Gencho’s generative formulation naturally supports soft acoustic matching involving different controls. Our contributions are as follows:

- We propose Gencho, a complex spectrogram-based diffusion-transformer that generates diverse, plausible RIRs from reverberant speech.
- We propose an enhanced design of a reverberation-structure-aware audio encoder that improves the model’s matching accuracy and generalization.
- We evaluate the designs and demonstrate Gencho’s flexible applications to acoustic matching, RIR completion with hybrid-prompting, and text-to-RIR generation.

Work done while Jackie Lin and Nishit Anand were in Adobe Internship.

¹Listening examples at <https://linjac.github.io/Gencho/>

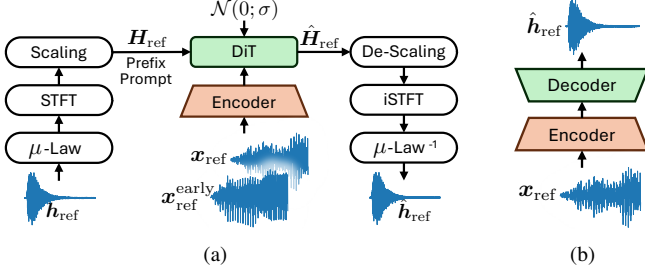


Fig. 1: (a) Model architecture of Gencho, the proposed generative estimator. (b) The non-generative FiNS-based baseline.

2. METHOD

2.1. Blind Room Impulse Response Estimation

We tackle the problem of blind RIR estimation, which takes a reference reverberant speech $x_{\text{ref}} = h_{\text{ref}} * s_{\text{ref}}$ and estimates the impulse response h_{ref} as a waveform signal without accessing the corresponding clean speech s_{ref} or prior information about the acoustic space. The estimated RIR is then convolved with the source content s_{source} to apply the acoustics acquired from the reference recording. Our blind RIR estimator aims to work on 48 kHz audio and assumes a 1.0 s duration for the impulse response. It consists of the structure-aware audio encoder and diffusion-based generative decoder. Figure 1 illustrates a schematic diagram of our generative approach and the non-generative model variants explored in this paper.

2.1.1. Reverberation Structure-Aware Audio Encoder

We first encode the reference recording from the target acoustic environment x_{ref} into a 128-dimensional latent embedding w_{ref} that captures the environment’s acoustic characteristics. To this end, we adopt the time-domain convolutional encoder design from FiNS[4] that consists of convolutional blocks, each consisting of a 1-D convolution with a kernel size of 15 and a stride of 2, and PReLU activation. It uses adaptive average pooling at the final output of the blocks to aggregate time frames into a global latent embedding. As a modification, we replace the original batch normalization at every block with a layer normalization to handle fatal failure cases, e.g., degenerate predictions when the input audio contains a significant amount of silence. Note, as shown in Fig. 1, that the diffusion-based model’s decoder is conditioned by the embedding vectors from the audio encoder also used in the non-generative method.

Our key design improvement to the audio encoder is separating the early reverberation component from the input audio. One challenge in blind estimation is to accurately model individual acoustic properties without mixing them up. An RIR signal can be decomposed into two parts, $h_{\text{ref}} = h_{\text{ref}}^{\text{early}} + h_{\text{ref}}^{\text{late}}$, which display drastically different structures in the early reflection component $h_{\text{ref}}^{\text{early}}$ (i.e., sparse) and the late tail component $h_{\text{ref}}^{\text{late}}$ (i.e., diffuse noise alike). Different kinds of nuanced errors, while displaying reasonable distance in the waveform or spectrogram space, could lead to drastically different perceptions by listeners. For example, early reflection errors contribute to the coloration and the sense of distance, while differences in reverb time are more noticeable when the tail is short.

Recent advances in audio separation and speech enhancement technologies enable fairly accurate isolation of foreground speech from background noise, with the flexibility to treat the perceptually sensible speech as foreground while leaving diffuse-noise-alike re-

verb tails to the background noise. Drawing on this capability and loosely inspired by the common practice of modeling the early reflection and the late tail of an IR separately [4, 6], we depart from prior methods that process the reverberant speech signal as a whole. Instead, we leverage a speech enhancement tool to extract the early-reflected component $x_{\text{ref}}^{\text{early}} = h_{\text{ref}}^{\text{early}} * s_{\text{ref}}$ from the reverberant speech x_{ref} . The extracted component is a relatively dry speech preserving the early reflection (e.g., the first 50ms of the reverberation) and coloration of the input. Based on this, the audio encoder takes in a two-channel input, comprising the extracted early component $x_{\text{ref}}^{\text{early}}$ and the full reverberant speech x_{ref} , and encodes it into a global embedding w_{ref} . This structured input enables the model to explicitly analyze different components of reverberation, resulting in more accurate and perceptually realistic room impulse response estimation.

2.1.2. Diffusion-based Generative Decoder

Diffusion-based models have drawn attention due to their generation capabilities in domains such as images [8] and audio [9, 10]. Diffusion models employ a forward process to add Gaussian noise to the data representation and learn a backward process to remove the noise via a neural network. In this work, we introduce a diffusion-based decoder that conditions on the global embedding w_{ref} from the audio encoder and generates the complex spectrogram of impulse responses using a diffusion process. The diffusion model enables modeling of the complex distribution of impulse response signals, producing various plausible yet perceptually consistent impulse responses. This approach captures the natural variability of room acoustics, avoiding the collapse behavior observed in the conventional regression-based methods.

Complex Spectrogram Representation. Since RIR’s spectrogram coefficients are sparse and of high variance, input normalization is essential: all one second-long target RIRs h_{ref} are μ -law encoded in the waveform domain to emphasize low amplitudes, transformed via short-time Fourier transform (STFT) with a frame size of 128 and a hop size of 64 samples to yield a complex spectrogram $H \in \mathbb{C}^{65 \times 751}$. Each spectrogram coefficient c is then power compressed $\tilde{c} = \beta |c|^\alpha e^{j\angle c}$ with empirically chosen $\alpha = 0.3$ and $\beta = 2$. The final input dimension is $\mathbb{R}^{130 \times 751}$ where the real and imaginary components are separated and stacked. We also experimented with a latent diffusion transformer using VAE latents learned from generic audio. However, it did not work as well as complex spectrogram diffusion, likely because the semantic relationships and distances in the general audio VAE space do not correlate strongly with perceptual distances in RIRs.

Diffusion Formulation. We use the v -prediction reparameterization [11] for the diffusion training objective, which has been shown to be more stable and effective than the naive reconstruction objective in producing high-quality outputs in audio generation domains [12]. The model structure is a Diffusion Transformer (DiT) [13] that offers significant advantages in scalability and robustness. The model conditions on the timestep embedding of the current diffusion step via adaptive layer normalization, and consists of transformer layers, each comprising RMS normalization, self-attention, cross-attention, another RMS normalization and a linear layer. It cross-attends to the encoded audio embedding from the encoder to guide the generation. We use classifier-free guidance during training, dropping out the condition with 10% probability, to encourage the model to learn both the conditional and unconditional distribution of RIRs.

Conditioning. We condition the diffusion transformer on the global embedding from the audio encoder via cross-attention. During train-

ing, we experimented with two initialization setups using the pre-trained audio encoder obtained from the non-generative FiNS-like model training: training from scratch and warm initialization. Warm initialization overall shows faster convergence and more stability.

RIR Completion and Hybrid Approach. During training of the diffusion model, we apply prefix prompting with a 50% probability, replacing up to the first 50ms of diffusion latent frames with the corresponding ground-truth RIR. This encourages the model not only to learn to generate an RIR from scratch, but also to "complete" an RIR from its early reflection component (i.e., as in [14]), offering versatile capabilities in downstream tasks. The model generates various versions of tails that sound perceptually coherent with the given partial RIR signal, allowing diverse data augmentation for acoustic simulation. We also develop a hybrid approach (see Sec. 4) that combines the benefits of the non-generative approach and the generative approach: we use the FiNS-like formulation to predict the early reflections while using the diffusion model to complete the tail.

3. EXPERIMENTS

3.1. Experiment Setup

The Baseline Model Variants. Three variants of FiNS were implemented as the regression-based baselines: (1) the original model, (2) a version with batch normalization layers replaced by layer normalization as described above, and (3) the layer normalization variant modified for two-channel input (the separated early reflected speech and the reverberant speech). Training was performed with a batch size of 256 for 60k steps using the AdamW optimizer with exponential decay scheduling starting at $lr = 0.0001$ with decay rate $\gamma = 0.999996$. We use the multiresolution STFT loss following the original FiNS.

The Diffusion Model Variants. The transformer-based backbone of the diffusion model has a hidden size of 256, comprised of 8 layers with 8 attention heads each, with learned positional embeddings, qk normalization, and a dropout probability of 0.1. Two diffusion model variants were implemented: (1) a single-channel model and (2) a dual-channel model. The diffusion encoders were warm-initialized with the trained encoder weights from FiNS variants (2) and (3), respectively, as they share the same encoder architecture. The diffusion model was trained for 100k steps with batch size = 256 on 8-gpu A100 using an AdamW optimizer and with linear cosine scheduling starting at $lr = 0.0001$ and decaying to 0.00001.

For inference, we adopt a DPM++ 2M SDE sampler with a Karras [15] noise schedule, using a log-SNR range of +5.0 to -8 and $\rho = 7.0$ to shape the noise decay, along with 24 sampling steps and a classifier-free guidance scale of 3.0.

3.2. Datasets

Training Datasets. The models were trained using a collection of room impulse responses from OpenSLR28 [16], the MIT IR Survey dataset [17], EchoThief [18], Arni [19], dEchorate [20], and the ACE Challenge dataset [21]. Since many of these datasets contain only a dozen unique rooms, but thousands of similar-sounding RIRs (mainly due to variations in microphone placement), we balanced the training data by weighting the contribution of each dataset according to its number of unique rooms, preventing overrepresentation of any single acoustic environment. Reverberant speech inputs were created by convolving these RIRs with clean speech from LibriTTS-R [22], a 585-hour high-quality dataset that was further upsampled

to 48 kHz using bandwidth extension [23]. To increase acoustic diversity, we applied data augmentations to RIRs and clean speech signals based on prior work [24, 25], randomly modifying the gains of direct-to-reverberant ratio and reverberation time durations for RIRs as well as scaling clean speech with varying speeds and volumes. All RIR signals were standardized to one second in length and normalized to a unit direct arrival energy at 2.5 ms, and the reverberant speech inputs consist of 5-second segments.

Evaluation Datasets. We evaluated the models on an out-of-domain dataset to examine their generalization capabilities. The test set consists of RIRs from BUTReverbDB [26] and OpenAIR [27] convolved with clean speech from the DAPS dataset [28], which in total span over 50 different environments, containing different signal characteristics and acoustic distributions from the training time.

4. RESULTS AND DISCUSSION

We report the standard RIR metrics of the methods. For each target-generated RIR pair, we compute the reverberation time to -60dB (T60), early decay time to -10dB (EDT), direct to reverberant ratio (DRR), and clarity index (C50). The percentage absolute error (PAE) of each time-related metric, and mean squared error (MSE) and mean absolute error (MAE) of every metric are reported in Tab. 1 for the out-of-domain test set. First, we see that the original FiNS does not perform well on new, unseen data; the layernorm and the structure-aware encoder improve performance, achieving 15.5% PAE on the EDT, 2.29 dB and 1.41 dB. Our diffusion-based model Gencho+2ch (taking in both the early reflected speech and the full reverberant speech) shows benefits in more accurately modeling the reverberation time, achieving the lowest T60 error with 13.6% PAE and 0.18s MAE, but degrades on EDT, DRR, and C50. We hypothesize that diffusion models are effective at modeling the randomness of noise-like tails in reverberation, but their inherent stochasticity may introduce excessive variability for early reflections, which tend to have clearer and more structured patterns.

Meanwhile, the RIR statistics on the evaluation dataset plotted in Fig. 2 show that our generative model is more expressive and approximates the underlying distribution of the target RIRs more closely than the non-generative baselines. The joint distribution of T60 and DRR reveals that FiNS exhibits a stronger inverse correlation between the two. This is likely a result of "regression to the mean", where the model tends to default to the common, averaged pattern

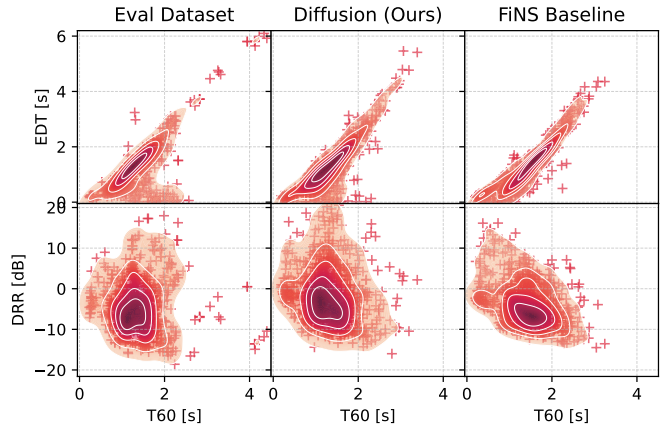


Fig. 2: Distribution of T60 vs EDT, and T60 vs DRR of the evaluation, our generated, and the FiNS layernorm generated samples.

Table 1: Comparison of the non-generative baseline and our diffusion model across T60, EDT, DRR, and C50 metrics.

Model	T60			EDT			DRR		C50	
	PAE (%)	MSE (s)	MAE (s)	PAE (%)	MSE (s)	MAE (s)	MSE (dB)	MAE (dB)	MSE (dB)	MAE (dB)
FiNS (original)	52.3	1.33	0.68	45.1	2.27	0.43	15.92	2.79	53.10	3.47
FiNS layernorm	22.6	0.19	0.31	20.1	0.14	0.22	8.93	2.38	3.69	1.45
FiNS layernorm + 2ch	14.2	0.10	0.20	15.5	0.08	0.16	8.86	2.29	3.89	1.41
Gencho + 1ch	16.3	0.13	0.23	26.4	0.22	0.28	25.91	4.11	12.15	2.69
Gencho + 2ch	13.6	0.08	0.18	34.8	0.16	0.25	25.62	4.04	11.81	2.66
Hybrid + 2ch @5ms	13.5	0.085	0.184	23.6	0.14	0.23	26.75	3.53	14.16	2.63
Hybrid + 2ch @25ms	13.8	0.085	0.185	25.5	0.13	0.23	23.78	3.40	12.09	2.45

that less reverberant environments are typically smaller and therefore associated with higher DRRs.

Hybrid approach. As observed, the non-generative baseline achieves more accurate DRR, whereas the diffusion-based approach produces more realistic decaying reverberation tails as reflected by improved T60 accuracy. Motivated by these complementary strengths, we construct a hybrid approach leveraging the IR completion capability of the trained diffusion model. We first estimate the RIR using the improved FiNS, then feed the early portion of that as a prompt to Gencho to generate the remaining RIR. We report the performance of this hybrid approach for prompt lengths = 5 ms and 25 ms in Tab. 1. The results show that this hybrid approach improves the accuracy of the generative model in EDT and DRR while maintaining T60 accuracy.

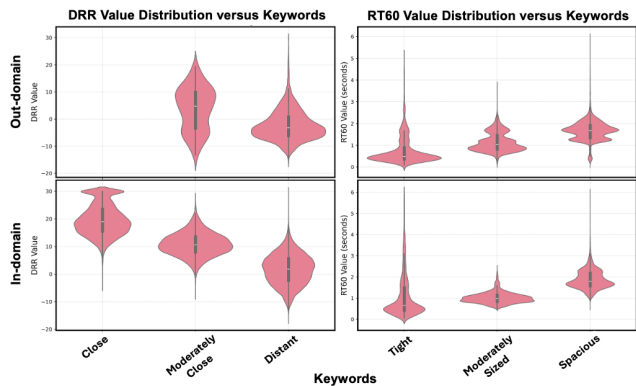
5. TEXT-TO-RIR GENERATION

Our diffusion-based RIR generator enables the creation of diverse, novel impulse responses under weak guidance and can be extended to a broad range of multi-modal applications. It offers flexible controls over acoustic properties and allows users greater freedom to explore virtual acoustic spaces. To demonstrate, we adapt the model for text-conditioned RIR generation (i.e., text-to-RIR) by replacing the audio encoder with a Flan-T5-XXL² text encoder. The diffusion model cross-attends to the text embedding sequence. This allows natural-language descriptions of an acoustic environment, e.g. “a large cathedral with long echoes and a distant human talker”, to condition the diffusion process to generate semantically matching RIRs.

The training dataset comprises approximately 150K pairs of impulse response signals and their corresponding text captions generated by ParaLLM [29] using a similar set of RIRs as in the aforementioned experiments. Each caption describes the semantic interpretation of the acoustic characteristics of the RIR in natural language, such as the perceived spaciousness of the room, the speaker’s apparent distance, and the clarity of the speech. We fine-tune a pre-trained Gencho decoder together with the newly introduced text encoder to enable transfer learning and accelerate convergence. We then evaluate the model on the validation RIR-caption pairs held-out from training as well as on a set of newly generated RIR-caption pairs covering unseen RIRs from OpenAIR [27].

Fig. 3 illustrates the strong correlations between reverberance-related keywords in the text descriptions and the acoustic properties (DRR and T60) of the generated RIRs across both test sets. Each keyword results in a differently centered distribution of acoustic properties that is semantically consistent with the keyword. For example, the keyword “tight (space)” leads to shorter reverberation

time, while “spacious” leads to longer reverb tails. OpenAIR does not contain closely captured impulse responses, so the keyword “close” does not appear in its captions. We further evaluate the accuracy based on the acoustic property bins designed in ParaLLM: **DRR Level** as DRR values divided into ≤ 5 , $5 \sim 11$, and ≥ 11 dB ranges, **T60 Level** as ≤ 0.5 , $0.5 \sim 1.2$, and ≥ 1.2 in seconds, and lastly a 3-by-3 grid of overall **Reverb Level** crossing DRR level and T60 level. We generated five variations of RIRs per text caption, achieving average accuracies of 70.1% and 89.1% on **DRR Level** for in-domain and out-domain test sets, 82.7% and 85.1% on **T60 Level**, and 57.2% and 76.3% on **Reverb Level**. The results show strong alignment with the ground-truth categories of the RIRs paired with the text captions, while the in-domain test set generally yields lower accuracies due to its broader thus more challenging coverage of acoustic environments. We refer readers to our demo page to check out the text-to-RIR examples.

**Fig. 3:** Violin plot. x-axis is short text prompts perceptually related to reverberance. y-axis is the DRR and T60 of generated RIRs.

6. CONCLUSION

We introduced Gencho, a diffusion transformer for blind room impulse response estimation that generates complex spectrogram RIRs from reverberant speech. The model leverages a structured input of separated early and late reflections to improve the robustness of conditioning information, and generates diverse, perceptually realistic outputs via a diffusion process. Gencho integrates seamlessly with standard speech-processing pipelines and can be adopted to a range of novel controllable tasks, including RIR completion and text-conditioned RIR generation. Experiments demonstrate improved generalization and richer acoustic statistics compared to non-generative baselines, establishing Gencho as a flexible tool for controllable acoustic simulation and generative audio applications.

²<https://huggingface.co/google/flan-t5-xxl>

7. REFERENCES

- [1] N. J. Bryan, “Impulse response data augmentation and deep neural networks for blind room acoustic parameter estimation,” in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 1–5.
- [2] P. Götz, C. Tuna, A. Brendel, A. Walther, and E. A. P. Habets, “Blind acoustic parameter estimation through task-agnostic embeddings using latent approximations,” in *2024 18th International Workshop on Acoustic Signal Enhancement (IWAENC)*, 2024, pp. 289–293.
- [3] C. Ick, A. Mehrabi, and W. Jin, “Blind acoustic room parameter estimation using phase features,” in *2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [4] C. J. Steinmetz, V. K. Ithapu, and P. Calamia, “Filtered noise shaping for time domain room impulse response estimation from reverberant speech,” in *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pp. 221–225, ISSN: 1947-1629.
- [5] A. Ratnarajah et al., “Towards improved room impulse response estimation for speech recognition,” in *2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [6] S. Lee, H.-S. Choi, and K. Lee, “Yet another generative model for room impulse response estimation,” in *2023 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2023, pp. 1–5.
- [7] A. Ratnarajah, S. Ghosh, S. Kumar, P. Chiniya, and D. Manocha, “AV-RIR: Audio-Visual Room Impulse Response Estimation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024, pp. 27164–27175.
- [8] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10684–10695.
- [9] Z. Evans et al., “Stable audio open,” *arXiv preprint arXiv:2407.14358*, 2024.
- [10] Z. Novack, J. McAuley, T. Berg-Kirkpatrick, and N. J. Bryan, “DITTO: Diffusion Inference-Time T-Optimization for Music Generation,” in *Forty-first International Conference on Machine Learning*, 2024.
- [11] T. Salimans and J. Ho, “Progressive distillation for fast sampling of diffusion models,” in *International Conference on Learning Representations*, 2022.
- [12] E. Hoogeboom, J. Heek, and T. Salimans, “Simple diffusion: End-to-end diffusion for high resolution images,” in *International Conference on Machine Learning*. PMLR, 2023, pp. 13213–13232.
- [13] W. Peebles and S. Xie, “Scalable diffusion models with transformers,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4195–4205.
- [14] J. Lin, G. Götz, and S. J. Schlecht, “Deep room impulse response completion,” *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2025, no. 1, pp. 20, 2025.
- [15] T. Karras, M. Aittala, T. Aila, and S. Laine, “Elucidating the design space of diffusion-based generative models,” *Advances in neural information processing systems*, vol. 35, pp. 26565–26577, 2022.
- [16] T. Ko, V. Peddinti, D. Povey, M. L. Seltzer, and S. Khudanpur, “A study on data augmentation of reverberant speech for robust speech recognition,” in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 5220–5224.
- [17] J. Traer and J. H. McDermott, “Statistics of natural reverberation enable perceptual separation of sound and space,” *Proceedings of the National Academy of Sciences*, vol. 113, no. 48, pp. E7856–E7865, 2016.
- [18] “EchoThief [dataset],” <https://www.echothief.com/echothief/>, [Accessed: 2024-09-29].
- [19] P. Karolina, S. J. Schlecht, and V. Välimäki, “Dataset of impulse responses from variable acoustics room aml at aalto acoustic labs,” Aug. 2022.
- [20] D. D. Carlo et al., “dEchorate: a calibrated room impulse response dataset for echo-aware signal processing,” *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2021, no. 1, pp. 39, 2021.
- [21] J. Eaton, N. D. Gaubitch, A. H. Moore, and P. A. Naylor, “The ace challenge — corpus description and performance evaluation,” in *2015 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2015, pp. 1–5.
- [22] Y. Koizumi et al., “LibriTTS-R: A Restored Multi-Speaker Text-to-Speech Corpus,” in *Interspeech*, 2023, pp. 5496–5500.
- [23] J. Su, Y. Wang, A. Finkelstein, and Z. Jin, “Bandwidth extension is all you need,” in *2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 696–700.
- [24] H. Yang, J. Su, M. Kim, and Z. Jin, “Genhancer: High-fidelity speech enhancement via generative modeling on discrete codec tokens,” in *Interspeech*, 2024.
- [25] J. Su, Z. Jin, and A. Finkelstein, “HiFi-GAN-2: Studio-quality speech enhancement via generative adversarial networks conditioned on acoustic features,” in *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2021, pp. 166–170.
- [26] I. Szöke, M. Skácel, L. Mošner, J. Paliesek, and J. Černocký, “Building and evaluation of a real room impulse response dataset,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 4, pp. 863–876, 2019.
- [27] S. Shelley and D. Murphy, “OpenAIR: An interactive auralization web resource and database,” *129th Audio Engineering Society Convention 2010*, pp. 1270–1278, 2010.
- [28] G. J. Mysore, “Can we automatically transform speech recorded on common consumer devices in real-world environments into professional production quality speech?—a dataset, insights, and challenges,” *IEEE Signal Processing Letters*, vol. 22, no. 8, pp. 1006–1010, 2014.
- [29] N. Anand et al., “Listening between the lines: Towards paralinguistic understanding of speech,” *2026 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2026.