

MB-RIRs: a Synthetic Room Impulse Response Dataset with Frequency-Dependent Absorption Coefficients

Enric Gusó,^{1,2} Joanna Luberadzka,² Umut Sayin,² Xavier Serra¹

¹ Universitat Pompeu Fabra, Music Technology Group, Barcelona
enric.guso@upf.edu, xavier.serra@upf.edu

² Eurecat, Centre Tecnològic de Catalunya, Tecnologies Multimèdia, Barcelona
joanna.luberadzka@eurecat.org, umut.sayin@eurecat.org

Abstract—We investigate the effects of four strategies for improving the ecological validity of synthetic room impulse response (RIR) datasets for monoaural Speech Enhancement (SE). We implement three features on top of the traditional image source method-based (ISM) shoebox RIRs: multiband absorption coefficients, source directivity and receiver directivity. We additionally consider mesh-based RIRs from the SoundSpaces dataset. We then train a DeepFilterNet3 model for each RIR dataset and evaluate the performance on a test set of real RIRs both objectively and subjectively. We find that RIRs which use frequency-dependent acoustic absorption coefficients (MB-RIRs) can obtain +0.51dB of SDR and a +8.9 MUSHRA score when evaluated on real RIRs. The MB-RIRs dataset is publicly available for free download.

1. INTRODUCTION

Speech enhancement (SE) is the combination of processes like dereverberation, denoising, declipping and bandwidth extension. In the last decade, deep learning-based SE methods have shown to be very effective, a success partially driven by the use of bigger and more diverse training datasets —e.g. the ones used in the Deep Noise Suppression Challenges (DNS) [1]. This data scaling has produced models that generalize better and mitigate cross-dataset performance differences like the ones in [2], [3]. On top of scaling size and variety, speech and noise datasets have also been extended to a sampling rate of 48kHz, further improving the results [4]–[6]. Generally speaking, the traditional approach for augmenting data is to convolve the speech signals with Room Impulse Responses (RIRs), simulating sounds in different acoustic environments.

However, most RIRs currently used in the state-of-the-art are still the ones from the early DNS challenges [7] —i.e. shoebox-like rooms with a single acoustic absorption coefficient for the entire frequency spectrum. They are typically rendered through the image source method and at 16kHz sampling rate, and were originally intended to train automatic speech recognition systems. In fact, increasing the realism of the RIRs has been shown to improve speech recognition for augmented reality glasses [8] or keyword spotting [9] even as a post-processing augmentation step [10], so the same could happen for the SE task.

More recently, the URGENT Challenge [6] acknowledges the importance of RIR generalization by using real recorded RIRs for evaluation. While it constitutes an important step for improving the ecological validity of SE models, it still uses DNS5 RIRs for training, so their models might be confined to that particular RIR coverage. The topic is also covered in [11], with a focus on augmenting existing RIRs with generative models.

Table 1: Summary of the evaluated synthetic RIR datasets. f_s stands for sampling rate, *rec* for receiver, *src* for source, *render* for rendering method and *T60* for reverberation time.

	f_s	<i>T60</i>	directivity		render
			<i>rec</i>	<i>src</i>	
DNS5	16k	single	✗	✗	ISM
SB	48k	single	✗	✗	ISM
MB	48k	multi	✗	✗	ISM
REC+MB	48k	multi	✓	✗	ISM
SRC+REC+MB	48k	multi	✓	✓	ISM
SSPA	48k	single	✓	✗	mesh

Beyond the monoaural SE task, mesh-based RIR data generation and refinement methods have been proposed in [12], [13], for tasks such as audio-visual navigation. An example of this is the SoundSpaces dataset [14], which contains a set of RIRs from rooms with complex geometries, providing a grid of positions for each room, and is rendered by path-based methods on 3D meshes. However, the performance of SE models trained on these RIRs has not been evaluated yet.

In this work we focus on the effects of RIRs used in training for SE. We evaluate five models trained on the same speech and noise datasets but on six different RIR datasets (two baselines, three of our own making and one from SoundSpaces), with the intention of answering whether models benefit from using RIRs with frequency dependent (multiband) absorption coefficients, whether modeling the receiver directivity or the source directivity helps to generalize in monoaural models [8], [9], and whether mesh-based modeling might be helpful or more suitable than Image Source Method (ISM) [9]. To this end, we propose an evaluation of the following RIR datasets, which are summarized above in Table 1.

- **DNS5:** state of the art (SOTA) baseline from [7].
- **SB:** we re-create DNS5 at 48kHz sampling rate using single-band absorption coefficients on the shoebox room material.
- **MB:** we use multi-band absorption coefficients instead.
- **REC+MB:** we add receiver directivity by using Head Related Transfer Functions (HRTFs) at the rendering stage.
- **SRC+REC+MB:** we add source directivity by modelling average human speech directivity in the ISM.
- **SSPA:** we use mesh-based RIRs from SoundSpaces instead.

2. TRAINING DATASETS

Regarding speech and noise, we have used the whole DNS5 which contains emotional speech, readings from English audiobooks using different accents, singing from VocalSet, French, Spanish, Italian, Russian and German speech from M-AILABS Speech, German speech from Wikipedia, Spanish from SLR73, SLR61, SLR39, SLR75, SLR74

The research leading to these results has received funding from the European union’s Horizon Europe programme under grant agreement No 101017884 - GuestXR project.

Dataset publicly available at <https://doi.org/10.5281/zenodo.15773093>

and SLR71 sets, which add up to a total of 583k speech utterances, and 63k noise samples from FreeSound and AudioSet [1] which add up to 1315 and 177 hours respectively. We have split them into 70% for training, and the remaining 30% split equally between validation and test sets, avoiding cross-set speaker contamination when possible. Below we describe the RIR training sets.

2.1. Baseline RIRs (DNS5)

We have taken as baseline 60k RIRs from DNS5 (synthetic RIRs from SLR26 and real RIRs from SLR28, originally at a 16kHz sampling rate and upsampled to 48kHz). We have distributed small, medium, and large rooms into three splits: 70%, 15% and 15% corresponding to training, validation and test sets.

2.2. Single-band absorption coefficients RIRs (SB)

With the intention of obtaining a fair comparison between the different strategies, we have created an approximation of the baseline DNS5 RIRs dataset using the Multichannel Acoustic Signal Processing (MASP) library [15], [16]. Specifically, we have generated 60k shoebox-like RIRs with single-band absorption coefficients, this time rendering at 48kHz sampling rate. The geometric configurations of these RIRs are described in Table 2: room dimensions \mathbf{r} have been sampled from uniform distributions, with r_x affecting r_y in order to avoid corridor-like geometries. Single-band absorption coefficients for all six walls of the shoebox-like room are computed via Sabine’s formula, using \mathbf{r} and reverberation time $T60$. In contrast to the multiband case (see 2.3), where $\mathbf{T60}$ is a vector defining reverberation time in each frequency band, in the single-band case $T60$ is a scalar computed as the mean of $\mathbf{T60}$ vector. $T60$ used for SB can be approximated by a normal distribution. Receiver coordinates \mathbf{rec} have also been set randomly for every room, avoiding to place them too close to the walls. The sources \mathbf{src} have been placed around the receivers at a distance between 0.5 and 3 meters and in front of them. We have used this same set of room configurations for the rest of the evaluated datasets except for the pre-computed SoundSpaces dataset SSPA.

2.3. Multiband absorption coefficients RIRs (MB)

The first of the three strategies we have proposed is to increase the dataset coverage by using multiband absorption coefficients $\mathbf{T60}$ [9] instead of a single $T60$ value as in the SB case. We have depicted the overall structure of the data generation pipeline in Figure 1. To ensure the generation of realistic parameters, we have analyzed 4495 real $\mathbf{T60}$ values from [17] in six frequency bands $\omega_n = \{125, 250, 500, 1k, 2k, 4k\}$ in Hz. We have modeled each band as an independent $Gamma(\alpha, \beta)$ distributions, fitted by minimizing the negative log-likelihood function. We have obtained shape $\alpha = \{1.72, 1.62, 1.93, 2.56, 4.17, 2.49\}$ and scale $\beta = \{0.39, 0.24, 0.14, 0.10, 0.09, 0.18\}$. MB RIRs are generated

Table 2: Random room configurations for SB, MB, REC+MB, SRC+REC+MB. Angles are in degrees, vectors are in bold.

$r_x = \mathcal{U}(3, 30)$	$rec_x = \mathcal{U}(0.35r_x, 0.65r_x)$
$r_y = r_x \cdot \mathcal{U}(0.5, 1)$	$rec_y = \mathcal{U}(0.35r_y, 0.65r_y)$
$r_z = \mathcal{U}(2.5, 5)$	$rec_z = \mathcal{U}(1, 2)$
$\ \mathbf{rec} - \mathbf{src}\ = \mathcal{U}(0.5, 3)$	$rec_\phi = \mathcal{U}(-45, 45)$
$\angle \mathbf{rec}, \mathbf{src} = \mathcal{U}(-45, 45)$	$rec_\theta = \mathcal{U}(-10, 10)$
$T60 \approx \mathcal{N}(0.4, 0.014)$	$T60 = \frac{1}{N} \sum_n \mathbf{T60}$
$\mathbf{T60} = Gamma(\alpha, \beta)$	

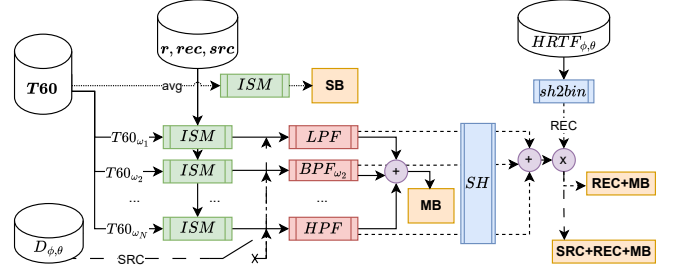


Fig. 1: RIR generation setup: the same set of geometric configurations is shared among the different datasets. Continuous line follows the MB pipeline, dotted line is for SB, short dashed line stands for the MB+REC modeling and long dashed line stands for the SRC modeling.

by running one ISM acoustic simulation for each sub-band, passing the sub-band RIRs through a filterbank comprised of a low-pass (LPF) for $n = 1$, band-pass (BPF_{ω_n}) for $n = \{2, \dots, N - 1\}$ and a high-pass filter (HPF) for the highest band, and finally summing the filtered RIRs.

2.4. Receiver directivity RIRs (REC+MB)

In [8], the augmented reality glasses’ directivity and the acoustical shadow of the head were shown to be relevant for speech separation and recognition. Likewise, SE from in-ear or headset devices may also benefit from modeling this directivity. Although there might be individual differences between the device types, they all share the characteristic of being placed near the human ear. We have modeled the receiver directivity in this scenario by applying a Head-Related Transfer Function (HRTF) to every reflection. To this end, we have used the same methodology as in [18]: to obtain the Spherical Harmonics (SH) expansion of the RIRs from MASP and then to use a Bilateral Magnitude Least Squares (BiMagLS) [19] decoder that implicitly applies the HRTFs, here focusing our evaluation to the left ear and using a set of normal hearing HRTFs of a Neumann KU100 dummy head.

2.5. Source directivity RIRs (SRC+REC+MB)

In SE, sources are always speech so we have applied its average directivity by taking the radiation pattern from [20] and converting it into the azimuth and elevation lookup table $D_{\phi, \theta}$. As depicted in Figure 1, $D_{\phi, \theta}$ can be applied in the SH pipeline right after the ISM method and prior to the filterbank rendering and summation. We have weighted the amplitude of each ISM reflection with the closest $D_{\phi, \theta}$. Both angles have been obtained through acoustical reciprocity [21]: swapping \mathbf{src} and \mathbf{rec} coordinates in the ISM we obtain a list of reciprocal reflections whose angle with respect to the receiver can be computed directly. Each reciprocal reflection angle with respect to the receiver is equivalent to the emission angle (with respect to the source) in the original ISM.

2.6. Mesh-based RIRs (SoundSpaces, SSPA)

The SSPA data set contains only 103 scenes, which is much less than scenes in SB. Nevertheless, we have included SSPA in this evaluation as a well-tested reference from fields like navigation [14] and 3D SE [22]. To our knowledge, it is the only publicly-available dataset which can be compared to REC+MB because it also uses Ambisonics for rendering into binaural and multiband absorption coefficients. It additionally provides complex geometries by path-tracing through 3D meshes. To keep the amount of RIRs consistent, we have taken 60k RIRs, also focusing on the left channel of the binaural downmix as in REC+MB and SRC+REC+MB.

3. EXPERIMENTAL SETUP

Given the speech, noise and the six RIR datasets described above, we have trained six DeepFilterNet3 [4] SE models – one for every RIR dataset– using its default hyperparameters except for a smaller maximum batch size of 38 in order to fit our GPUs and a p_reverb of 1 (to apply a RIR to every utterance). We have chosen DeepFilterNet3 because it has SOTA performance while being open source and having real-time capabilities. After training for 117 to 120 epochs (depending on the early stopping) our evaluation has been three-fold: firstly, we have applied the models to monoaural noisy utterances convolved with real RIRs and computed a set of intrusive and non-intrusive metrics. Secondly, we have used a few processed samples to evaluate subjectively with a MUSHRA listening test [23]. Thirdly, we have applied the models to the *headset* and *speakerphone* DNS5 test sets to address the effects of directivity. Since these test sets contain real mixtures (no separate clean speech ground truth is available) metrics are restricted to be non-intrusive.

3.1. Objective evaluation

Speech and noise for these simulations has been taken from the DNS5 test set splits. More precisely, we have taken 10k high quality read speech samples from the VCTK [27] and LibriVox [28] subsets (7k and 3k respectively). We have used 397 real RIRs from the ACE [29], MIT IR Survey [30], Openair [31], BUT Reverb [32] datasets and the RIRs from SLR28 [7] that were not used during training. Noises and RIRs have been reused as many times as necessary to fit the speech test set size.

From each speech sample we have taken a non-silent four-seconds chunk y , convolved it with a real RIR h and summed it to a noise sample n using $SNR = \mathcal{U}(0, 30)$, obtaining noisy and reverberant x . We have made sure that time synchronicity between x and y is kept, but when addressing this we have found that certain real RIR onsets are harder to detect than their non-noisy synthetic counterparts, so the vicinity of the highest amplitude has been checked for any earlier peak that surpassed a -6dB threshold. We have also found that the broadly used correlation-based synchronization method is less robust to RIR noise than our thresholding heuristic.

Note that in the traditional signal model $x = (y * h) + n$, noise can be non-reverberant while the speech is, which can constitute a limitation despite being broadly used in most of SE literature. To further investigate this, we have also conducted the evaluation using $x = (y + n) * h$, a signal model in which both speech and noise have a more spectrally-coherent reverberation at the expense of using the exact same RIR (and therefore placing speech and noise at the exact same position in the simulation space, which is deemed unfeasible). Results on this alternative signal model are not reported here for the sake of brevity, but were found to be very similar to the results from the more traditional $x = (y * h) + n$ we report below.

Regarding metrics, we have used almost all of the ones from the URGENT Challenge [6] with the addition of Scale Invariant Signal-to-Distortion Ratio ($SISDR$) and $SISDR_{squim}$, $PESQ_{squim}$, MOS_{squim} and $STOI_{squim}$ from [24]. Because reverberation can mislead some metrics' values and our dataset is the first dataset of this kind that we are aware of, we are reporting the results as increments with respect to the noisy and reverberant signal. Given DeepFilterNet3 model f and clean speech estimate $\hat{y} = f(x)$, we report evaluation metrics m increments as $\Delta m = m(\hat{y}) - m(x)$ for non-intrusive metrics and $\Delta m = (m(\hat{y}, y) - m(x, y))$ for the intrusive ones that require the reference y . For DeepFilterNet3 absolute metrics on the commonly-used Voicebank+Demand test set we refer the reader to [4].

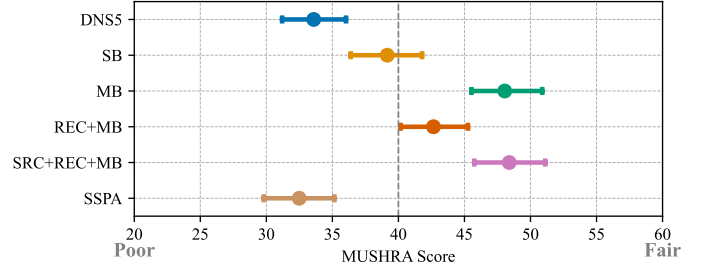


Fig. 2: MUSHRA scores mean and 95% confidence intervals for each model.

3.2. Subjective evaluation

In order to further validate the objective evaluation, we have picked eight reverberant synthetic speech and noise mixtures that use real RIRs from the same set as in the objective evaluation and have conducted an online MUSHRA listening test using [23], [33], taking y as 48kHz high-quality clean reference and hidden reference, x as anchor and the six different model estimates \hat{y} using a scale from 0 or Bad to 100 or Excellent. All 33 participants have declared to be expert listeners (professional audio producers or researchers with experience doing listening tests). Only one subject had to be excluded during the post-screening due to misjudging the anchor.

3.3. Computational requirements

Due to the sequential nature of the ISM and the computational demands of high-order SH, generating all the RIR datasets simultaneously has required a server with 64-cores of CPU and 250GB of RAM for 10 days. Training all the six DeepFilterNet3 models took 30 days on two A100s.

4. RESULTS AND DISCUSSION

Subjective results of the enhanced speech and noise mixtures convolved with real RIRs are depicted on Figure 2.

To compare the conditions, we used a pairwise t-test with a significance threshold of $p < 0.05$. MUSHRA scores show no statistically significant difference between DNS5 and SSPA (t-test: $p = 0.55$) which points out that variety and quantity of rooms are factors to take into account on top of RIR complexity –note that SSPA additionally models complex room geometries and furniture but despite having the same number of RIRs than SB, it contains much fewer different rooms. SB significantly outperforms the DNS5 baseline (t-test: $p = 0.003$), which highlights the importance of extending the sampling rate from 16kHz to 48kHz.

All models apart from SSPA have received significantly better scores than DNS5 baseline (SB: $p = 3 \cdot 10^{-3}$, MB: $p = 7 \cdot 10^{-14}$, REC+MB $p = 8 \cdot 10^{-7}$ and SRC+REC+MB $p = 2 \cdot 10^{-14}$), indicating that both the higher sampling rate and adding MB, SRC and REC features into the acoustic simulation improves the speech enhancement. Interestingly, there are no significant differences between MB and SRC+REC+MB (t-test: $p = 0.86$) which suggests that modeling directivity does not degrade monoaural SE performance when evaluated on real RIRs. We can't assess SRC and REC directivities with the real RIR test set, but note that MB and SRC+REC+MB outperform the rest and that all multiband RIR datasets perform similarly. The scores are presented in Table 3, in addition to objective intrusive and non-intrusive results.

Similar to the MUSHRA scores, the objective results on the same set (the real RIRs test set) indicate that MB RIRs outperform the rest. However, there is no consensus among all intrusive metrics, nor between intrusive and non-intrusive metrics. While most intrusive

Table 3: Evaluation results on real RIRs using DNS5 read speech. Mean values \pm standard deviation. The higher the better (\uparrow) except for Log-Spectral Distance and Mel-Cepstral Distortion (\downarrow). SQUIM [24], NISQA [25] and DNSMOS [26] are non-intrusive neural-based approximations of the metrics.

		DNS5	SB	MB	REC+MB	SRC+REC+MB	SSPA
Intrusive							
$\Delta SISR$	\uparrow	1.105 \pm 3.50	1.335 \pm 3.49	1.361 \pm 3.45	1.965 \pm 4.52	1.9 \pm 4.41	0.888 \pm 3.27
ΔSDR	\uparrow	2.817 \pm 3.17	3.298 \pm 3.02	3.328 \pm 3.00	2.674 \pm 3.19	2.706 \pm 3.17	2.025 \pm 3.11
ΔLSD	\downarrow	1.213 \pm 1.41	1.34 \pm 1.42	1.26 \pm 1.40	0.841 \pm 1.43	0.778 \pm 1.43	1.105 \pm 1.43
ΔMCD	\downarrow	-1.434 \pm 2.00	-1.479 \pm 1.98	-1.504 \pm 1.99	-1.542 \pm 1.97	-1.533 \pm 1.98	-1.454 \pm 2.00
$\Delta PESQ$	\uparrow	0.52 \pm 0.47	0.583 \pm 0.47	0.593 \pm 0.47	0.595 \pm 0.46	0.609 \pm 0.46	0.51 \pm 0.45
$\Delta STOI$	\uparrow	0.088 \pm 0.07	0.101 \pm 0.07	0.102 \pm 0.07	0.101 \pm 0.06	0.101 \pm 0.06	0.082 \pm 0.06
$\Delta PhonSim$	\uparrow	0.051 \pm 0.15	0.058 \pm 0.15	0.058 \pm 0.15	0.063 \pm 0.15	0.0626 \pm 0.15	0.061 \pm 0.15
$\Delta SpkSim$	\uparrow	-0.055 \pm 0.13	-0.058 \pm 0.14	-0.057 \pm 0.13	-0.034 \pm 0.12	-0.03 \pm 0.11	-0.053 \pm 0.13
$\Delta BertSim$	\uparrow	0.081 \pm 0.07	0.086 \pm 0.07	0.087 \pm 0.07	0.089 \pm 0.07	0.09 \pm 0.07	0.082 \pm 0.07
Non-intrusive							
$\Delta SISR_{squim}$	\uparrow	3.793 \pm 5.87	4.456 \pm 5.71	4.544 \pm 5.67	3.814 \pm 5.65	4.038 \pm 5.63	2.291 \pm 5.61
ΔMOS_{squim}	\uparrow	0.547 \pm 0.71	0.537 \pm 0.71	0.541 \pm 0.71	0.523 \pm 0.70	0.506 \pm 0.71	0.541 \pm 0.70
ΔMOS_{dnsmos}	\uparrow	0.884 \pm 0.53	0.935 \pm 0.54	0.937 \pm 0.53	0.895 \pm 0.53	0.906 \pm 0.53	0.811 \pm 0.52
ΔMOS_{nisqa}	\uparrow	1.158 \pm 0.73	1.198 \pm 0.73	1.224 \pm 0.72	1.109 \pm 0.72	1.099 \pm 0.72	0.977 \pm 0.73
$\Delta PESQ_{squim}$	\uparrow	0.636 \pm 0.62	0.667 \pm 0.59	0.678 \pm 0.59	0.613 \pm 0.58	0.637 \pm 0.58	0.525 \pm 0.57
$\Delta STOI_{squim}$	\uparrow	0.06 \pm 0.08	0.07 \pm 0.08	0.071 \pm 0.08	0.063 \pm 0.08	0.066 \pm 0.08	0.048 \pm 0.08
Subjective							
<i>MUSHRA</i>	\uparrow	33.59 \pm 20.3	39.15 \pm 22.6	48.05 \pm 22.8	42.65 \pm 22.5	48.39 \pm 23.0	32.49 \pm 21.7

metrics follow the subjective scores and show a higher performance of MB and SRC+REC+MB datasets, *SISR*, *MCD* and Phonetic Similarity (*PhonSim*) metrics suggest to use REC+MB. Non-intrusive metrics agreement is higher than in the intrusive ones, with MB slightly outperforming the rest (t-test: $p < 0.05$), perhaps due to a bias in NISQA and SQUIM training data but interestingly, towards MB instead of SB.

Comprehensive results for the DNS5 *headset* and *speakerphone* are shown in Table 4. A priori, one would expect SRC+REC+MB to outperform the rest for *headset*, but this is not the case. For both *headset* and *speakerphone*, most non-intrusive metrics show MB models slightly outperforming the rest as they did for the real RIR test set, but differences are not statistically significant (e.g. t-test: $p = 0.9$ for *SISR*_{squim} between MB and SB). Therefore, we can not find evidence of benefits from directivity modeling.

Table 4: Non-intrusive evaluation results on real noisy and reverberant recordings from the DNS5 *headset* and *speakerphone* sets. *HDS* stands for *headset* and *SPK* for *speakerphone*. The higher the better for all metrics.

		DNS5	SB	MB	REC+MB	SRC+REC+MB	SSPA
<i>SISR</i> _{squim}	<i>HDS</i>	15.65	15.75	15.80	15.35	15.42	14.21
	<i>SPK</i>	16.81	16.86	17.05	16.65	16.82	15.51
<i>MOS</i> _{squim}	<i>HDS</i>	3.92	3.89	3.90	3.91	3.94	3.91
	<i>SPK</i>	4.15	4.15	4.15	4.14	4.14	4.14
<i>MOS</i> _{dnsmos}	<i>HDS</i>	3.01	3.04	3.05	3.04	3.02	3.01
	<i>SPK</i>	3.03	3.05	3.06	3.03	3.04	2.99
<i>MOS</i> _{nisqa}	<i>HDS</i>	3.54	3.64	3.68	3.61	3.56	3.43
	<i>SPK</i>	3.77	3.82	3.82	3.77	3.76	3.62
<i>PESQ</i> _{squim}	<i>HDS</i>	2.58	2.59	2.62	2.55	2.54	2.45
	<i>SPK</i>	2.64	2.63	2.64	2.56	2.59	2.49
<i>STOI</i> _{squim}	<i>HDS</i>	0.92	0.92	0.92	0.92	0.92	0.91
	<i>SPK</i>	0.94	0.94	0.94	0.94	0.94	0.93

Although in [8] they showed that directivity could be exploited even in monoaural models, we could not replicate the results, either because of the limitations of NISQA and SQUIM metrics, or due to limitations in the DNS5 test set. In both cases, the lack of a clear improvement when applying directivity brings into question if MB RIRs perform better thanks to being more domain consistent or just because they are more diverse [34] than SB. We would like to address this in future work, perhaps by comparing acoustically implausible random *T60* RIRs with the ones from Section 2.3.

Overall, we recommend to train SE models with high sampling rate multiband RIRs, potentially incorporating SRC and REC directivities depending on the use case. We make both MB and SRC+REC+MB available at <https://doi.org/10.5281/zenodo.15773093>.

5. CONCLUSIONS

We have presented an evaluation of three Room Impulse Response (RIRs) generation techniques using the state of the art real-time capable Speech Enhancement (SE) model DeepFilterNet3. Specifically, we have shown that the idea of extending the RIR coverage to frequency dependent acoustic absorption coefficients –which has been shown to be successful for keyword spotting– can also benefit modern SE models. We have also found that the amount and variability of RIRs can be as important as the model complexity, showing that a lot of ISM-based rooms can outperform fewer mesh-based rooms (SoundSpaces) when evaluated on real RIRs. Besides, we have shown that source and receiver directivities do not degrade performance in the monoaural SE task and also have the potential of increasing SE performance but could not provide strong evidence of the improvement. In conclusion, we recommend training with a varied 48kHz multiband RIRs dataset: MB-RIRs.

REFERENCES

- [1] H. Dubey, A. Aazami, V. Gopal, B. Naderi, S. Braun, R. Cutler, A. Ju, M. Zohourian, M. Tang, M. Golestaneh *et al.*, “ICASSP 2023 deep noise suppression challenge,” *IEEE Open Journal of Signal Processing*, vol. 142, pp. 725–737, Mar. 2024.
- [2] J. Richter, S. Welker, J.-M. Lemercier, B. Lay, and T. Gerkmann, “Speech enhancement and dereverberation with diffusion-based generative models,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 31, pp. 2351–2364, Jun. 2023.
- [3] B. Kadioğlu, M. Horgan, X. Liu, J. Pons, D. Darcy, and V. Kumar, “An empirical study of Conv-TasNet,” in *Proc. ICASSP*, May 2020, pp. 7264–7268.
- [4] H. Schröter, T. Rosenkranz, A. Maier *et al.*, “DeepFilterNet: Perceptually motivated real-time speech enhancements,” in *Proc. Interspeech*, Aug. 2023, pp. 2008–2009.
- [5] W. Zhang, K. Saijo, Z.-Q. Wang, S. Watanabe, and Y. Qian, “Toward universal speech enhancement for diverse input conditions,” in *Proc. ASRU*, Dec. 2023, pp. 1–6.
- [6] W. Zhang, R. Scheibler, K. Saijo, S. Cornell, C. Li, Z. Ni, J. Pirklbauer, M. Sach, S. Watanabe, T. Fingscheidt, and Y. Qian, “URGENT challenge: Universality, robustness, and generalizability for speech enhancement,” in *Proc. Interspeech*, Sep. 2024, pp. 4868–4872.
- [7] 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 *Proc. ICASSP*, Mar. 2017, pp. 5220–5224.
- [8] R. Arakawa, M. Parvaix, C. Lai, H. Erdogan, and A. Olwal, “Quantifying the effect of simulator-based data augmentation for speech recognition on augmented reality glasses,” in *Proc. ICASSP*, Apr. 2024, pp. 726–730.
- [9] E. Bezzam, R. Scheibler, C. Cadoux, and T. Gisselbrecht, “A study on more realistic room simulation for far-field keyword spotting,” in *Proc. APSIPA ASC*, Dec. 2020, pp. 674–680.
- [10] A. Ratnarajah, Z. Tang, and D. Manocha, “TS-RIR: Translated synthetic room impulse responses for speech augmentation,” in *Proc. ASRU*, Dec. 2021, pp. 259–266.
- [11] J. Lin, G. Götz, H. S. Llopis, H. Hafsteinsson, S. Guðjónsson, D. G. Nielsen, F. Pind, P. Smaragdis, D. Manocha, J. Hershey *et al.*, “Generative data augmentation challenge: Synthesis of room acoustics for speaker distance estimation,” in *Proc. ICASSPW*, Apr. 2025, pp. 1–5.
- [12] Z. Tang, R. Aralikatti, A. J. Ratnarajah, and D. Manocha, “Gwa: A large high-quality acoustic dataset for audio processing,” in *Proc. SIGGRAPH*, Aug. 2022, pp. 1–9.
- [13] L. Kelley, D. Di Carlo, A. A. Nugraha, M. Fontaine, Y. Bando, and K. Yoshii, “RIR-in-a-box: Estimating room acoustics from 3D mesh data through shoebox approximation,” in *Proc. Interspeech*, Sep. 2024, pp. 3255–3259.
- [14] C. Chen, C. Schissler, S. Garg, P. Kobernik, A. Clegg, P. Calamia, D. Batra, P. Robinson, and K. Grauman, “Soundspaces 2.0: A simulation platform for visual-acoustic learning,” *Proc. NeurIPS*, vol. 35, pp. 8896–8911, Nov. 2022.
- [15] A. Perez-Lopez and A. Politis, “A python library for multichannel acoustic signal processing,” in *Audio Eng. Soc. Conv. 148*, May 2028.
- [16] A. Politis, *Microphone array processing for parametric spatial audio techniques*. Espoo, Finland: Aalto University, 2016.
- [17] K. Sridhar, R. Cutler, A. Saabas, T. Parnamaa, M. Loide, H. Gamper, S. Braun, R. Aichner, and S. Srinivasan, “ICASSP 2021 acoustic echo cancellation challenge: Datasets, testing framework, and results,” in *Proc. ICASSP*, Jun. 2021, pp. 151–155.
- [18] E. Gusó, J. Luberadzka, M. Baig, U. Sayin, and X. Serra, “An objective evaluation of hearing aids and dnn-based binaural speech enhancement in complex acoustic scenes,” in *Proc. WASPAA*, Oct. 2023, pp. 1–5.
- [19] I. Engel, D. Goodman, and L. Picinali, “Improving binaural rendering with bilateral ambisonics and MagLS,” in *Annual German Conference on Acoustics*, vol. 99, 2021, p. 10.
- [20] T. W. Leishman, S. D. Bellows, C. M. Pincock, and J. K. Whiting, “High-resolution spherical directivity of live speech from a multiple-capture transfer function method,” *J. Acoust. Soc. Am.*, vol. 149, pp. 1507–1523, Mar. 2021.
- [21] P. Samarasinghe, T. D. Abhayapala, and W. Kellermann, “Acoustic reciprocity: An extension to spherical harmonics domain,” *J. Acoust. Soc. Am.*, vol. 142, no. 4, pp. EL337–EL343, Oct. 2017.
- [22] C. Marinoni, R. F. Gramaccioni, C. Chen, A. Uncini, and D. Comminiello, “Overview of the L3DAS23 challenge on audio-visual extended reality,” in *Proc. ICASSP*, Jun. 2023, pp. 1–2.
- [23] B. Series, “Method for the subjective assessment of intermediate quality level of audio systems,” *International Telecommunication Union Radiocommunication Assembly*, vol. 2, 2014.
- [24] A. Kumar, K. Tan, Z. Ni, P. Manocha, X. Zhang, E. Henderson, and B. Xu, “Torchaudio-squim: Reference-less speech quality and intelligibility measures in torchaudio,” in *Proc. ICASSP*, May 2023, pp. 1–5.
- [25] G. Mittag, B. Naderi, A. Chehadi, and S. Möller, “NISQA: A deep cnn-self-attention model for multidimensional speech quality prediction with crowdsourced datasets,” in *Proc. Interspeech*, Aug. 2021, pp. 2958–1796.
- [26] C. K. Reddy, V. Gopal, and R. Cutler, “DNSMOS: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors,” in *Proc. ICASSP*, Jun. 2021, pp. 6493–6497.
- [27] C. Veaux, J. Yamagishi, K. MacDonald *et al.*, “CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit,” *University of Edinburgh. The Centre for Speech Technology Research*, vol. 6, p. 15, Nov. 2017.
- [28] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an ASR corpus based on public domain audio books,” in *Proc. ICASSP*, Apr. 2015, pp. 5206–5210.
- [29] J. Eaton, N. D. Gaubitch, A. H. Moore, and P. A. Naylor, “The ACE challenge — corpus description and performance evaluation,” in *Proc. WASPAA*, Oct. 2015, pp. 1–5.
- [30] I. Szöke, M. Skácel, L. Mošner, J. Paliesek, and J. Černocký, “Building and evaluation of a real room impulse response dataset,” *IEEE J. Sel. Top. Signal Process.*, vol. 13, pp. 863–876, Aug. 2019.
- [31] S. Shelley and D. T. Murphy, “Openair: An interactive auralization web resource and database,” in *129th Audio Eng. Soc. Convention*, vol. II, Nov. 2010, pp. 1270–1278.
- [32] I. Szöke, M. Skácel, L. Mošner, J. Paliesek, and J. Černocký, “Building and evaluation of a real room impulse response dataset,” *IEEE J. Sel. Top. Signal Process.*, vol. 13, no. 4, pp. 863–876, May 2019.
- [33] D. Barry, Q. Zhang, P. W. Sun, and A. Hines, “Go listen: an end-to-end online listening test platform,” *Journal of Open Research Software*, vol. 9, p. 20, Jul. 2021.
- [34] C.-B. Jeon, G. Wichern, F. G. Germain, and J. Le Roux, “Why does music source separation benefit from cacophony?” in *Proc. ICASSPW*, Apr. 2024, pp. 873–877.