

Mean absorption estimation from room impulse responses using virtually supervised learning^{a)}

Cédric Foy,¹ Antoine Deleforge,^{2,b)} and Diego Di Carlo³

¹Unité mixte de recherche en acoustique environnementale (UMR AE), Cerema, Université Gustave Eiffel, Institut français des sciences et technologies des transports (Ifsttar), Strasbourg, 67035, France

²Université de Lorraine, Centre national de la recherche scientifique (CNRS), Inria, Laboratoire lorrain de recherche en informatique et ses applications (LORIA), Nancy, F-54000, France

³Université Rennes, Inria, CNRS, Institut de recherche en informatique et systèmes aléatoires (IRISA), France

ABSTRACT:

In the context of building acoustics and the acoustic diagnosis of an existing room, it introduces and investigates a new approach to estimate the mean absorption coefficients solely from a room impulse response (RIR). This inverse problem is tackled via *virtually supervised learning*, namely, the RIR-to-absorption mapping is implicitly learned by regression on a simulated dataset using artificial neural networks. Simple models based on well-understood architectures are the focus of this work. The critical choices of geometric, acoustic, and simulation parameters, which are used to train the models, are extensively discussed and studied while keeping in mind the conditions that are representative of the field of building acoustics. Estimation errors from the learned neural models are compared to those obtained with classical formulas that require knowledge of the room's geometry and reverberation times. Extensive comparisons made on a variety of simulated test sets highlight different conditions under which the learned models can overcome the well-known limitations of the diffuse sound field hypothesis underlying these formulas. Results obtained on real RIRs measured in an acoustically configurable room show that at 1 kHz and above, the proposed approach performs comparably to classical models when reverberation times can be reliably estimated and continues to work even when they cannot. © 2021 Acoustical Society of America. <https://doi.org/10.1121/10.0005888>

(Received 1 February 2021; revised 28 July 2021; accepted 28 July 2021; published online 20 August 2021)

[Editor: Peter Gerstoft]

Pages: 1286–1299

I. INTRODUCTION

When sound propagates in a room, its reflections on the walls, ceiling, floor, and other surfaces lead to the well known phenomenon of reverberation. When the reverberation level is too high, it can be a major source of nuisance for the room's users. To alleviate this, some of the main parameters an acoustician can act on are the *absorption coefficients* of the room surfaces, namely, the proportion of sound energy that the surfaces' materials do not reflect. These are generally frequency dependent and typically expressed within octave bands, $b \in \mathcal{F} = \{0.125, 0.25, 0.5, 1, 2, 4\}$ kHz, in room acoustics standards. To obtain the acoustic diagnosis of a room and deduce a renovation plan, acousticians need to know the absorption coefficients $\alpha_i(b)$ of each individual surface i in the room. This is typically done through a manual iterative process in which acoustic simulators are tuned to match the *in situ* measurements while taking into account the room's geometry and properties of known materials as measured in laboratories.

Among the *in situ* measurements used in practice, room impulse responses (RIR) are rich signals that capture the

acoustic signature of the room via the shape of their decay, their echo density over time, or the timings of their early echoes. Although the *forward* physical process from the acoustic parameters to RIRs is well understood, as illustrated by the existence of many reasonably accurate and efficient RIR simulators (Habets, 2006; Scheibler *et al.*, 2018; Schimmel *et al.*, 2009), the *inverse problem* of retrieving the absorption coefficients of surfaces solely from a RIR is much more challenging and is the focus of this article. We consider the simple but common case of a *shoobox* (cuboid) room with a different material on each of the six surfaces. Even in this case, recovering the absorption coefficients of all surfaces from a single RIR without any knowledge on the source, receiver, or wall positions is out of reach because of inherent ambiguities of the problem such as permutations between the different surfaces. To alleviate this issue, this work focuses on estimating the area-weighted mean absorption coefficients,

$$\bar{\alpha}(b) = \frac{\sum_i \alpha_i(b) S_i}{\sum_i S_i} \in [0, 1], \quad (1)$$

where S_i denotes the area of surface i in m^2 . Note that this quantity is treated, here, as a purely analytical parameter that globally summarizes the acoustic properties of all surfaces in

^{a)}This paper is part of a special issue on Machine Learning in Acoustics.

^{b)}Electronic mail: antoine.deleforge@inria.fr, ORCID: 0000-0003-0339-7472.

the room. In acoustics, it is traditionally used under the hypothesis of a *diffuse sound field* (DSF) in which the energy is uniformly distributed in space and flows isotropically (Kuttruff, 2009; Nolan *et al.*, 2018). However, in this work, we will also consider its estimation under more general, non-diffuse settings. Choosing this particular quantity as a target will notably allow relevant comparisons to methods based on the classical reverberation theory, i.e., by inverting the well-known Sabine and Eyring formulas (Kuttruff, 2009), at least under conditions that are close to the DSF regime.

We propose to tackle the inverse problem of estimating $\bar{\alpha} = [\bar{\alpha}(b)]_{b \in \mathcal{F}} \in [0, 1]^6$ from a single RIR without any other information on the room using supervised machine learning and, in particular, nonlinear regression.

Whereas artificial neural networks have proven to be a very powerful family of models for nonlinear regression in recent years, a well-known bottleneck is their need for a large number of input-output pairs to be trained. As of today, because *in situ* estimation of absorption coefficients remains a costly and complex task, sufficiently large and diverse real RIR databases annotated with surface absorption profiles are not available. Hence, we propose to make use of *virtually supervised learning* as introduced in Gaultier *et al.* (2017). The idea is to use the known forward physical model, namely, a room acoustic simulator, to generate a potentially unlimited amount of annotated data to learn the inverse mapping from. The main contributions of this article are (i) a novel approach to efficiently sample simulated training data that are representative of commonly encountered acoustics in cuboid rooms, which is shown to outperform naive uniform sampling; (ii) an extensive comparative simulation study between estimates based on classical reverberation theory and those obtained from various neural network designs, including their generalizability to unseen data, noise, and various acoustic conditions; and (iii) a comparative study between virtually trained models and classical models on real measured RIRs.

Our simulated experiments reveal that neural models can successfully estimate the mean absorption coefficients under a wide range of acoustical conditions with the mean absolute errors below 0.05 while not requiring any geometrical information on the room. As expected, in non-DSF settings, they are more accurate than classical models that rely on the DSF hypothesis. On real data that are close to the DSF regime, errors obtained from the proposed learned model are not satisfying below 1 kHz but remain under 0.1 in higher octave bands and are comparable to those obtained with classical models. Moreover, in those higher frequencies, it is shown that the neural model continues to yield reliable $\bar{\alpha}(b)$ estimates even in conditions where classical models cannot as reverberation times cannot be extracted from RIRs due to the lack of sufficient linear decays in the Schroeder curves (Schroeder, 1965).

Although the observed limitations of classical formulas from reverberation theory outside of the DSF regime are well-known and expected (Nolan *et al.*, 2018), they still constitute an interesting comparison point as these tools

remain widely used today to obtain initial *in situ* acoustical estimates in practice, e.g., Prawda *et al.* (2020). Further investigation on the real-world applicability of learned models in lower octave bands and their extension to the geometrically informed estimation of individual absorption profiles are left for future work.

The remainder of this work is organized as follows. Section II provides an overview of the related works. Section III details the construction of our simulated RIR datasets, examining trade-offs between computational tractability, realism, and representativity. Section IV presents the neural networks' design and training. Sections V and VI contains our extensive comparative experimental study on both the simulated and real data. Finally, Sec. VII concludes and offers leads for future work.

II. RELATED WORKS

A. Absorption coefficient estimation

Whereas this article focuses on the intermediate task of estimating the area-weighted mean absorption coefficients in a room, the estimation of the individual absorption coefficients or, more generally, the surface impedance of a material is a vast and long-standing research topic, which is briefly reviewed here. The most commonly used techniques require an isolated sample of the studied material in a controlled environment. The impedance tube method is one of the most widely used methods (ASTM, 2006; ISO, 2001), and the associated analytical approach is usually that of Chung and Blaser (1980a,b), based on the transfer function between two microphones. Alternatively, the *reverberation room* method (ISO, 2003) uses the theory of reverberation and relies on the DSF hypothesis.

In contrast, this article explores *in situ* estimation. For a recent exhaustive review of this topic, the reader is referred to Brandão *et al.* (2015). Classically, the goal is to separate the direct wave from the reflected wave in an impulse response with different constraints that depend on the acoustic environment. Early approaches include *echo-impulse* methods in which the reflected wave is extracted by eliminating the incident wave and parasite wave using temporal windowing or subtraction. Due to the time-frequency uncertainty relation $\Delta t \Delta f \geq 1$ (Garai, 1993), a compromise must then be found between the size of the time-domain filters used and the information loss at low frequencies. Also, to have a good temporal separation of the waves, the emitted pulse must be narrow, of flat frequency spectrum, and repeatable, which is difficult to have in practice (Cramond and Don, 1984; Davies and Mulholland, 1979; Garai, 1993; Yuzawa, 1975).

To overcome these limitations, methods based on stationary noise have been proposed. Although Barry (1974) and Hollin and Jones (1977) use white noise, Aoshima (1981) and Suzuki *et al.* (1995) later proposed a flat spectrum pulse signal stretched in time by filtering. Other excitation signals were then developed to guarantee a better immunity to background noise such as maximum-length sequences (MLS) (Rife and Vanderkooy, 1999; Schroeder,

1979; Stan *et al.*, 2002) and sine sweep signals (Farina, 2000, 2007; Müller and Massarani, 2001). To date, the advantages and disadvantages of these signals are still being studied (Guidorzia *et al.*, 2015; Torras-Rosell and Jacobsen, 2010).

In parallel, other works focus on the development of analytical models of propagation. In Ingård and Bolt (1951), the sound field of an anechoic room is approximated by a set of plane waves. This was later reiterated in Ando (1968) and Sides and Mulholland (1971). Allard and Sieben (1985) introduced the microphonic doublet approach and specific impedance, which can be related to *surface impedance* using the linearized Euler equation. This approach is only valid if the distance between the microphones is small compared to the wavelength (Allard and Aknine, 1985; Champoux and L'espérance, 1988; Champoux *et al.*, 1988; Minton *et al.*, 1988). More finely, the sound field can be modeled by a set of spherical waves, as proposed in Champoux *et al.* (1988), based on the analytical model of Nobile and Hayek (1985) and later in Li and Hodgso, (1997). Finally, approaches based on the principle of acoustical holography, following Tamura (1990), have also been recently investigated (Nolan, 2020; Rathsam and Rafaely, 2015; Richard *et al.*, 2017). Although simple propagation models are easily invertible, more realistic propagation models are generally not, requiring the use of more complex and approximate numerical solvers, as well as access to precise details on the acoustic environment, which are not always available to field acousticians, in practice (Brandão *et al.*, 2015).

In summary, estimating the absorption coefficients of a material remains a complex task. It hinges on the choice of a number of parameters that are often correlated with each other and hard to precisely control, in practice, such as the excitation signal, source and receiver properties, environment (free field, anechoic, reverberant), experimental setup (number and position of sources and microphones, size of the material under study), chosen propagation model, and post processing. Developing a generic approach to retrieve absorption profiles *in situ* from a unique RIR measurement at an arbitrary location is, hence, an attractive research avenue for building acoustics.

B. Machine learning in acoustics

Machine learning methodologies have only recently been applied to acoustics. They are still relatively scarce in the field but have received fast growing interest (Bianco *et al.*, 2019). Whereas the lack of a large amount of training data is often a limiting factor, this has been alleviated by the use of massive simulations (Gaultier *et al.*, 2017; Kim *et al.*, 2017), data augmentation (Gamper and Tashev, 2018), or domain adaptation (He *et al.*, 2019). Early successful applications of machine learning to acoustics were mostly in sound source localization (Chakrabarty and Habets, 2017; Deleforge *et al.*, 2015a; Deleforge *et al.*, 2015b; Di Carlo *et al.*, 2019; Gaultier *et al.*, 2017; He *et al.*, 2019; Lefort *et al.*, 2017; Niu *et al.*, 2017) and acoustic scene and event classification (Deecke and Janik, 2006; Gradišek *et al.*,

2017; Mesaros *et al.*, 2017; Mesaros *et al.*, 2019; Parsons and Jones, 2000). The concept of *acoustic space learning* was introduced in Deleforge *et al.* (2015a) in the context of sound source localization. A large dataset of broadband audio recordings from different (source, receiver) locations in a fixed room was gathered using a motorized binaural head. A supervised nonlinear regression model was then trained on this dataset to learn a mapping from audio features to source directions. This approach is, however, limited by the data availability and does not generalize well to different acoustic environments as shown in Deleforge *et al.* (2015b). To alleviate this issue, the concept was later extended to *virtual acoustic space learning* (Gaultier *et al.*, 2017; Kataria *et al.*, 2017) in which hundreds of thousands of examples are generated using a room acoustic simulator. In the context of sound localization, such virtually learned models showed some direct, albeit limited, generalizability to real data in Gaultier *et al.* (2017) and Chakrabarty and Habets (2017). In He *et al.* (2019), a domain adaptation technique was proposed to strengthen this generalizability.

Closer to our application, supervised learning was recently used to estimate the reverberation time (Gamper and Tashev, 2018) or volume (Genovese *et al.*, 2019) of a room *blindly*, i.e., from the single channel noisy recording of an unknown speech source. Interestingly, these works use a careful combination of real and simulated data for training. Performances are, however, naturally limited in such blind settings. In a preliminary study (Kataria *et al.*, 2017), virtually supervised learning was used to jointly estimate the mean absorption coefficients of the walls and three-dimensional (3D) position of a broadband noise source from binaural recordings. The room shape, receiver position, and properties of the floor and ceiling were fixed and known throughout, whereas the absorption coefficients of the walls were supposedly frequency independent and only results on the simulated data were reported. Even more recently, a method to estimate the six absorption coefficients of the surfaces of a shoebox room in increasing order in a fixed frequency band from an impulse response was proposed using a fully connected deep neural network (Yu and Kleijn, 2021). The model was both trained and tested on simulated RIR datasets using the image source method without diffusion or noise and with absorption coefficients uniformly drawn at random between zero and one. Such an absorption distribution is, however, not representative of commonly encountered room acoustics as will be showed in Sec. III B. Reported errors were 30%–60% lower than random guessing, but no comparison to known acoustical models and no experiments on real data were performed.

III. SIMULATED DATASETS

The first step of the proposed virtually supervised approach is to simulate a large number of RIRs paired with corresponding mean absorption coefficients $\bar{\alpha}$ [Eq. (1)] to train our models. For this, two important trade-offs must to be considered. The first one is between the realism of simulations and their computational demand and is governed by

the choice of a simulator and tuning of its internal parameters. The second trade-off is between the diversity of considered acoustic environments and amount of representative data needed to train the model. Both trade-offs are discussed in detail in Secs. III A and III B.

A. Realism trade-off

When simulating RIRs, more realism typically implies higher, sometimes prohibitive computational costs. Existing room acoustic simulators can be divided into three categories (Habets, 2006). The first category solves the wave equation in discretized space, time, and/or frequency domains. These notably include finite element methods (Okuzono *et al.*, 2014), boundary-element methods (Pietrzyk, 1998), or finite-difference time-domain methods (Botteldooren, 1995). Whereas they can, in principle, simulate any acoustic conditions and geometry to arbitrary precision, their computational time depends on the space discretization steps used, which conditions the attainable wavelengths. In the context of building acoustics, which deals with frequencies as high as 5 kHz within large volumes, accurately generating thousands of RIRs is unfeasible with such methods. A second category includes variants of the well-known image source model, originally proposed in Allen and Berkley (1979), many times extended, e.g., Borish, (1984), Peterson (1986), Samarasinghe *et al.* (2018), and implemented in many widely used acoustic simulators, e.g., Habets (2006), Scheibler *et al.* (2018), and Schimmel *et al.* (2009). This deterministic method allows very efficient implementations, in particular, in cuboid rooms but only models ideal specular reflections on surfaces and, hence, lacks realism. The last category includes energetic methods based on Monte Carlo sampling, which is also known as ray-tracing or particle filtering (Kulowski, 1985; Schimmel *et al.*, 2009; Schröder, 2011). Similar to wave-based methods, these approaches can, in principle, model arbitrary acoustic conditions and are particularly well-suited to model surface scattering. However, their computational time and precision depends on the number of rays (or equivalently particles). For such methods to be tractable in the context of room acoustics, the receiver must typically be approximated by a large receptive field to aggregate enough rays. Alternatively, the diffuse-rain method proposed in Schröder (2011) systematically sends a proportion of diffuse energy to a point receiver at each ray collision, reducing the number of rays needed. In both cases, the timings of rays reaching the receivers are nondeterministic and only reflect acoustical effects in a statistical, energetic sense.

For this study, we choose a hybrid simulator belonging to the last two categories, which is referred to as Roomsim and proposed in Schimmel *et al.* (2009). Roomsim combines the image source method to obtain precise timings of specular reflections dominating the early part of the RIR and the diffuse-rain method to account for stochastic diffuse effects dominating the RIR's tail. The hybrid simulator Roomsim enables frequency-dependent absorption and scattering

coefficients and uses a minimum-phase finite-impulse-response representation of rays reaching the receiver to convert echograms into RIRs. This minimum-phase representation is physically motivated by the causality and fast-decaying properties of the resulting signals. A software based on Roomsim is shown to yield remarkably accurate RIRs compared to measured RIRs in identical conditions in Wabnitz *et al.* (2010). We used the open-source C++/MATLAB implementation from the original authors (Schimmel *et al.*, 2009). As a compromise between accuracy and computational demand, we used a frequency of sampling of 48 kHz, 50 000 rays per simulation for the diffuse-rain method, and image sources up to order 50 for the image source method. Simulations were run and aggregated along the following six octave bands: $b \in \mathcal{F}$. These match those available in most absorption coefficient databases and are commonly used in acoustic regulations. Although its impact is minor, atmospheric attenuation is taken into account for a temperature of 20 °C and a relative humidity of 42% (Roomsim default values).

We must stress that although lower frequencies are perceptually relevant in building acoustics, the energy-based simulation approach used here is unable to accurately model some of the wave phenomena occurring below Schroeder's frequency (Schroeder, 1996) such as room modes (Schröder, 2011, Sec. 5.6). This limitation of the current study will be reflected in our real-data experiments as discussed in Sec. VI.

B. Representativity trade-off

A large diversity in the training data is generally desirable to learn a model that generalizes well to many different situations. However, more diversity also implies more data to obtain a representative training dataset. Indeed, for a fixed sampling density of a parameterized observation space, the number of required samples grows exponentially in the number of parameters, an effect known as the *curse of dimensionality*. As a mitigating trade-off, we choose in this study to focus on environments that are representative of the field of building acoustics, e.g., offices, schools, restaurants, or accommodations. In particular, we exclude very large volumes such as those encountered in churches, tunnels, hangars, or swimming pools. Our evaluation will also exclude unusual absorption profiles that are only encountered in highly specialized rooms (e.g., anechoic or semi-anechoic chambers). Figure 1 shows the absorption profiles of the 92 commonly encountered reflective, wall, floor, and ceiling materials that will be considered in this study.¹ Because most commonly encountered rooms in buildings are cuboids, this study focuses on those rather than dealing with arbitrary complex geometries. This is also motivated by the fact that the image source method is much faster in this setting as exploited by Roomsim. Finally, we only consider empty rooms. This strong assumption is partially mitigated by the use of the diffuse-rain model. The random sound rays stemming from this Monte Carlo approach can

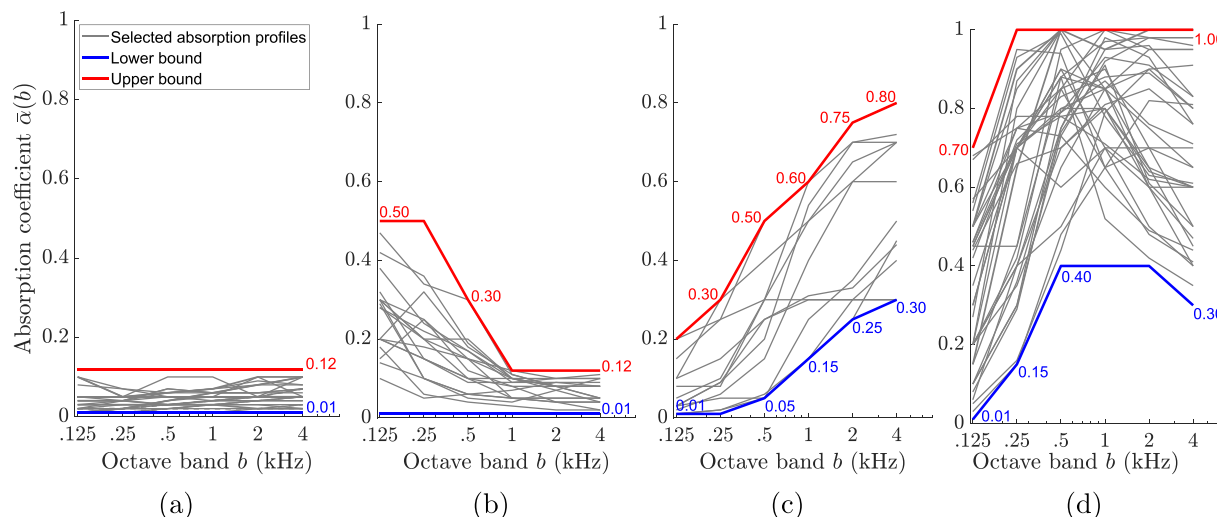


FIG. 1. (Color online) The absorption profiles of 92 commonly encountered reflective, wall, floor, and ceiling materials with lower and upper bounds. (a) 26 reflective profiles, (b) 19 wall profiles, (c) 12 floor profiles, and (d) 35 ceiling profiles are shown.

approximate reflections on objects of different sizes, depending on the octave bands/wavelengths considered.

The relevant parameters impacting RIRs can then be divided into a reasonably small set of geometric and acoustic parameters. The geometric parameters include the 3D positions of the source and receiver (both assumed to be omnidirectional in this study) and the width L_x , length L_y , and height L_z of the room. The height L_z was drawn uniformly at random between 2.5 m and 4 m, and the width L_x and length L_y were drawn uniformly at random between 1.5 m and 10 m. The receiver and source positions were drawn uniformly at random in the room for each RIR while ensuring a minimum distance of 0.5 m to any surface and 1 m between the two using rejection sampling (ISO, 2008).

The acoustic parameters include the absorption $\alpha_i(b)$ and scattering $s_i(b)$ coefficients of each of the six surfaces i in each of the six octave bands b . Two different strategies were explored to sample the absorption coefficients. The first, most straightforward strategy is to draw all 36 coefficients uniformly at random between 0 and 1 for each RIR. We later refer to this approach as *Unif*, which is also the approach employed in the recent paper by Yu and Kleijn (2021). The obtained $\bar{\alpha}(b)$ distribution [Eq. (1)] over 15 000 simulated RIRs is shown in Fig. 2(a). As can be observed in Fig. 2(b), the resulting histogram of $RT_{30}(b)$ values² is narrowly spread around 150 ms, which is an unusual value mostly encountered in semi-anechoic chambers. This is because using this technique, drawing four or more reflective absorption profiles within the same room [e.g., $\bar{\alpha}_i(b) < 0.15$ for all b] is very unlikely. Yet, highly reflective profiles are frequently encountered in real buildings. These are characteristics of hard surfaces made of, e.g., concrete, bricks, or tiles. The absorption profiles of 26 such materials are plotted in Fig. 1(a). As can be seen, they are all roughly frequency independent with absorption coefficients below 0.12. Based on this, we designed the following new *reflectivity biased* (RB) sampling strategy:

- (1) For each surface type (wall, floor, ceiling), toss a coin;
- (2) on heads, draw reflective frequency-independent absorption profiles uniformly at random in $[0.01, 0.12]$ for these surfaces; and
- (3) on tails, draw nonreflective frequency-dependent absorption profiles uniformly at random within predefined ranges, depending on the surface type (see Fig. 1).

Note that walls are either all reflective or all nonreflective but may still have distinct profiles. The nonreflective ranges are chosen to encompass typical materials used on walls, floors, and ceilings in common buildings as shown in Figs. 1(b)–1(d). As can be seen in Figs. 2(d) and 2(c), the proposed RB sampling technique results in more diverse and more representative distributions for both reverberation times $RT_{30}(b)$ and mean absorption coefficients $\bar{\alpha}(b)$. The peak around 0.06 observed in Fig. 2(c) is consistent with the proposed bias toward reflective surfaces and the chosen realistic absorption ranges.

Finally, for both the Unif and RB sampling strategies, the same frequency-dependent scattering profile was used for all surfaces. This approach, previously used in Gaultier *et al.* (2017), is based on the interpretation that the diffuse-rain model of Roomsim globally captures the random reflections in the room rather than the specific local effects. Whereas random scattering coefficients in $[0, 1]$ were used in all octave bands for Unif, we, respectively, used the ranges $[0, 0.3]$ and $[0.2, 1]$ for octave bands in $\{125 \text{ Hz}, 250 \text{ Hz}, 500 \text{ Hz}\}$ and $\{1 \text{ kHz}, 2 \text{ kHz}, 4 \text{ kHz}\}$ for RB. This choice is guided by scattering profiles measured in real rooms as reported in Vorländer and Mommertz (2000). Overall, 1 training set of 15 000 RIRs and 1 development set of 5000 RIRs were generated for each of the 2 sampling techniques.

IV. NEURAL NETWORK MODELS AND TRAINING

A. Data preprocessing

A crucial problem in supervised learning is that of finding an appropriate representation for input data, which is

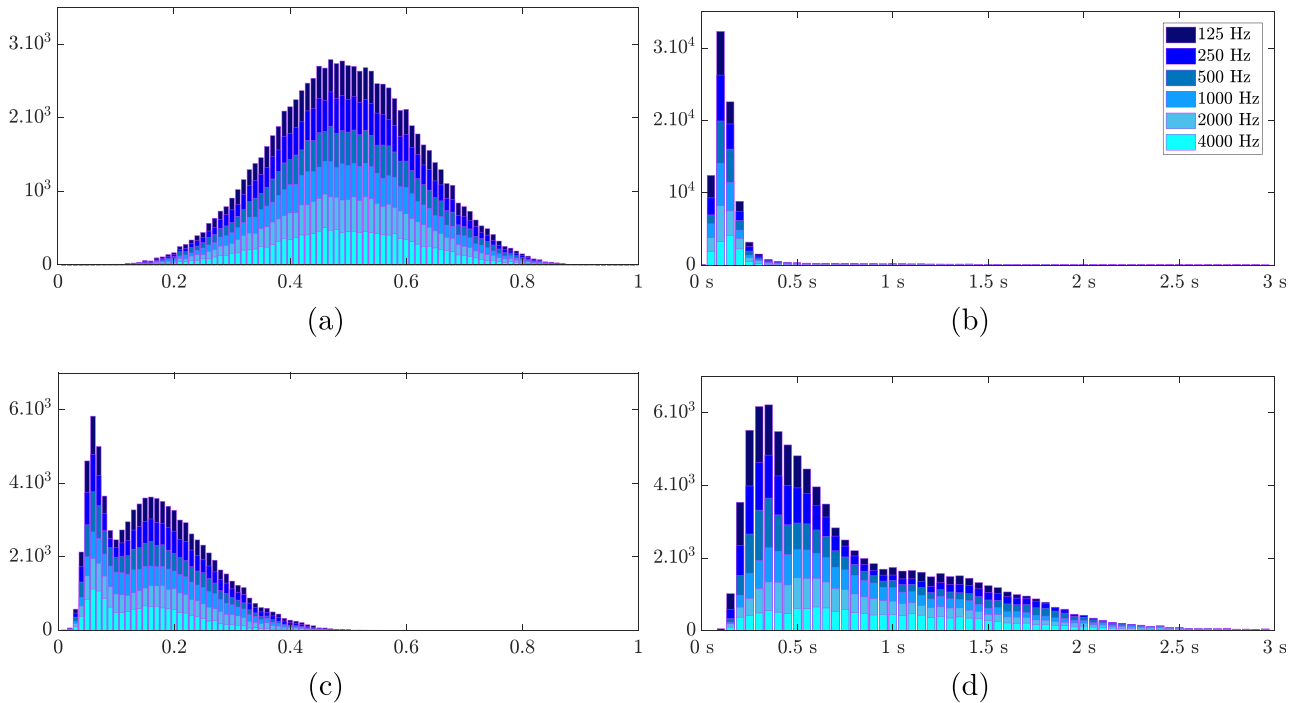


FIG. 2. (Color online) The histograms of $\bar{\alpha}(b)$ [(a),(c)] and $RT_{30}(b)$ [(b),(d)] values in 6 octave bands for 15 000 RIRs using Unif [(a),(b)] vs RB [(c),(d)] sampling.

sometimes referred to as the *feature extraction* step. Ideally, one seeks a representation that preserves or enhances features that are relevant for estimating the output while removing unnecessary or redundant features. In learning-based audio signal processing applications, phaseless time-frequency representations, such as magnitude spectrograms or Mel-frequency cepstral coefficients, have been widely used. Because frequency-dependent values are sought, such representations seem attractive at first glance. However, by discarding the phase, they would remove fine-grain temporal information such as the timings of early echoes in the RIRs. These timings could be exploited to infer geometrical properties of the room that, in turn, correlate with absorption coefficients conditionally on the reverberation time, as shown by Eq. (2). Alternatively, one could consider invertible complex time-frequency representations such as the short-term Fourier transform (STFT). Our preliminary experiments in that direction were, however, not conclusive, possibly the result of the difficulty of handling nonlinear complex phase behavior in the networks, or because any choice of STFT parameters implies a nonobvious compromise between the time and frequency resolution at each frame. Consequently, we choose to let the network learn its own internal representation of time-domain RIRs in an end-to-end fashion. This approach has recently shown considerable success in other audio signal processing applications, e.g., Luo and Mesgarani (2018).

RIRs obtained by Roomsim were resampled from 48 to 16 kHz. In fact, the highest octave band considered does not exceed 5.7 kHz, suggesting that 12 kHz could be sufficient for our application. However, higher-frequency features, such as the times of arrival of early reflections, may still

carry useful information. On the other hand, overly relying on very high frequencies would be disconnected from real applications as the receivers and emitters used to measure RIRs are always band limited in practice. Only the first 500 ms of the RIRs were preserved as this range is expected to contain the most salient acoustical information, including both early and late reflections. This resulted in 8000-dimensional input vectors. A random white Gaussian noise with a signal-to-noise ratio (SNR) of 30 dB was also added to every RIR in the datasets. This is expected to make learned models more robust and prevent them from relying on vanishingly small values in the RIRs, which would be inaccessible in practical applications. Finally, all input vectors were normalized to have a maximum value of one. This is done to facilitate learning and also prevent models from relying on the RIR's absolute amplitude, which is often inaccessible in practical applications due to unknown source and microphone gains.

B. Network design

Two commonly used neural network architectures are considered for this study, namely, the multilayer perceptron (MLP) and convolutional neural network (CNN), depicted in Figs. 3(a) and 3(b), respectively. The MLP is made up of three fully connected hidden layers of successive dimensions 128, 64, and 32, each followed by exponential linear units (ELUs). The CNN starts with three consecutive one-dimensional (1D)-convolutional hidden layers with a stride of 1, respective filter sizes 33, 17, 98 with zero-padding to preserve dimensionality after each convolution, and number of filters 64, 32, and 16. Each convolution is followed by a

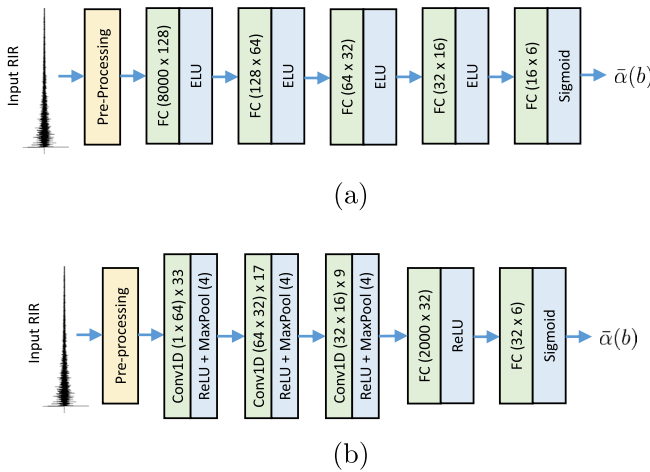


FIG. 3. (Color online) The neural network architectures. (a) The multilayer perceptron (MLP) and (b) convolutional neural network (CNN) are shown.

max pooling layer of width four and ELUs. The resulting output of dimension 2000 is then passed through a fully connected hidden layer of size 32 with ELUs. This particular design of layers is meant to define two simple dimensionality-reducing networks of relatively small and comparable sizes and depths. For each network, a final fully connected output layer is used to yield the desired output vector, evaluated by a mean-squared error loss function. Networks are optimized on the training set using batches of size of 1000 and ADAM (Kingma and Ba, 2014) with a learning rate of 0.001. Parameters, yielding the lowest average loss on the development set over 400 epochs, are used in all experiments. These meta-parameters and choice of ELUs, rather than rectified linear units (ReLUs), were guided by preliminary experiments on the development sets.

Three different output targets were considered: (i) the six-dimensional vector of mean absorption coefficients in all octave bands $\bar{\alpha} \in [0, 1]^6$, (ii) the vector of inverse mean absorption coefficient $\bar{\alpha}^{-1} \in \mathbb{R}^{+6}$, or (iii) the concatenation of the mean absorption and scattering coefficients $[\bar{\alpha}; \bar{s}] \in [0, 1]^{12}$. The second idea derives from the fact that the reverberation of a room is roughly inversely proportional to the mean absorption in DSF conditions, e.g., Sabine’s law (Kuttruff, 2009). The third idea is to test whether annotating the network with scattering coefficients at train time could help the estimation of absorption, i.e., multi-task learning. Output values in $[0, 1]$ were constrained using sigmoid gates, whereas positive values were constrained using a ReLU. A comparison of the distribution of absolute errors on $\bar{\alpha}$ obtained on the development set of RB using these three targets is shown in Fig. 4. In the remainder of this article, the absolute error is defined as the absolute difference between the target and estimated values. For a given dataset, reported means or box plots are computed over all input RIRs but also over all six octave bands unless stated otherwise. As can be seen, using inverse or concatenated vectors yields equivalent or worse results than simply using $\bar{\alpha}$. Hence, only networks outputting $\bar{\alpha}$ are considered in the remainder of this article.

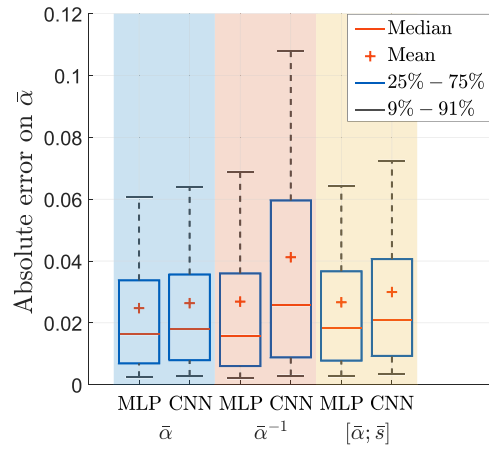


FIG. 4. (Color online) The comparison of three output layers, which were trained on RB and evaluated on the RB development set.

Figures 5(a) and 5(b) show the evolution of the loss functions of the two networks on the training and development sets for both Unif and RB. It can be observed that the MLP is more prone to over-fitting than the CNN. This suggests that the latter generalizes better to the unseen RIRs, an effect which will be confirmed in Sec. VI. This might be explained by the use of temporal convolutions, which may more efficiently capture the global frequency content of the RIRs than fully connected layers, while discarding less relevant local information.

V. EXPERIMENTS AND RESULTS

A. Baseline models

As a comparison point with the proposed neural models, we use mean absorption estimates obtained using the well-known Sabine’s law and its more precise variant from Eyring from the reverberation theory (Kuttruff, 2009):

$$\bar{\alpha}_{\text{Sabine}}(b) = 0.163V/(S), \tag{2}$$

$$\bar{\alpha}_{\text{Eyring}}(b) = -\ln(1 - \bar{\alpha}_{\text{Sabine}}(b)), \tag{3}$$

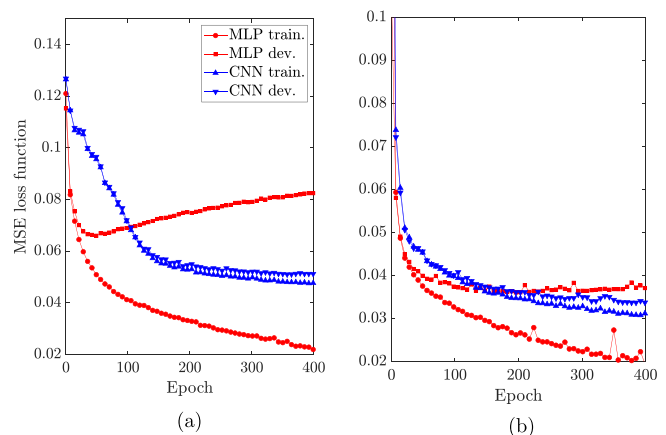


FIG. 5. (Color online) The loss evolution on training and development sets. (a) The Unif datasets and (b) RB datasets are shown.

where V denotes the room's volume and $S = \sum_i S_i$ is its total surface. Eyring's and Sabine's models are always given the true volume V and total surface S of the room in all experiments. Obviously, the DSF hypothesis inherent to these classical models is not theoretically verified for many of the considered room configurations. To better understand the impact of this limitation, a preliminary study was performed on the Unif and RB training databases. The reverberation time used in the formulas was calculated on different dynamics $([-5 \text{ dB}, -15 \text{ dB}], [-5 \text{ dB}, -20 \text{ dB}], [-5 \text{ dB}, -25 \text{ dB}], [-5 \text{ dB}, -35 \text{ dB}], [-5 \text{ dB}, -65 \text{ dB}])$ of the Schroeder curves (Schroeder, 1965), and the resulting distributions of the absolute errors were estimated. The dynamic $[-5 \text{ dB}, -35 \text{ dB}]$, i.e., $RT_{30}(b)$, was retained for our study as it offered the smallest median values of absolute errors for the Unif and RB training databases, i.e., 0.07 and 0.03, respectively. Such low errors show that the exploitation of these DSF-based models in our comparative study, while limited, is not unreasonable for the selected room configurations.

B. Simulation results

We now compare the different learned models (MLP-Unif, MLP-RB, CNN-Unif, CNN-RB) to Eyring's [Eq. (3)] and Sabine's [Eq. (2)] models on the task of estimating surface-weighted mean absorption coefficients [Eq. (1)] from a simulated RIR. A variety of simulated test sets,

containing 500 RIRs each and all generated with Roomsim, are considered.

The first simulated test set, called *realistic*, only contains surface materials commonly encountered in real buildings and is drawn uniformly at random from the database presented in Fig. 1. Five fixed geometries, representative of typical rooms, were selected for this set with the following (L_x, L_y, L_z) dimensions in meters: (4,5,3), (10,2,3), (10,5,3), (5,8,2.5), (10,10,5). The scattering of the walls and noise level are the same as those in the RB datasets. Absolute errors obtained with the six methods are presented in the form of box plots in Fig. 6(a). As can be seen, networks trained on the naive Unif training set do not succeed in outperforming the classical approaches based on reverberation theory. However, mean estimation errors, twice as small as Eyring's method and with much less variance, are obtained using the networks trained on the RB set. As expected, the estimates by Sabine are shown to be slightly less accurate than those of Eyring. Hence, results from the Unif-trained networks and Sabine's model will no longer be reported in what follows. The absolute error distribution was also observed per octave band for this test set (Fig. 7). No major differences in errors were observed across the octave bands for the different methods. Therefore, the errors will systematically be aggregated over all octave bands in the remainder of this section.

We then conduct a series of experiments on specially crafted simulated test sets to further assess the efficiency of

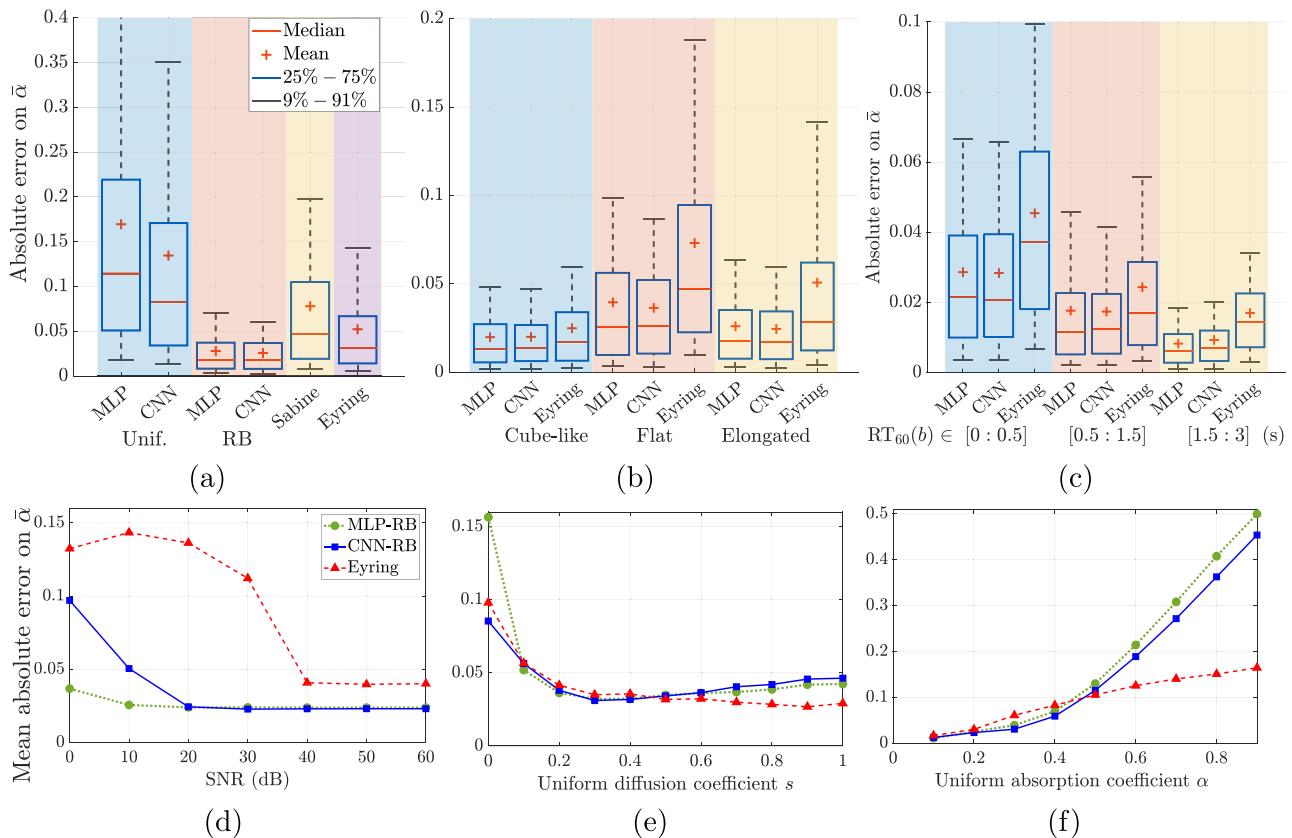


FIG. 6. (Color online) The comparison of $\bar{\alpha}$ estimation errors on different simulated test sets of 500 RIRs each. (a) The realistic test set, (b) influence of geometry, (c) influence of reverberation time, (d) influence of noise, (e) influence of diffusion, and (f) influence of absorption are shown.

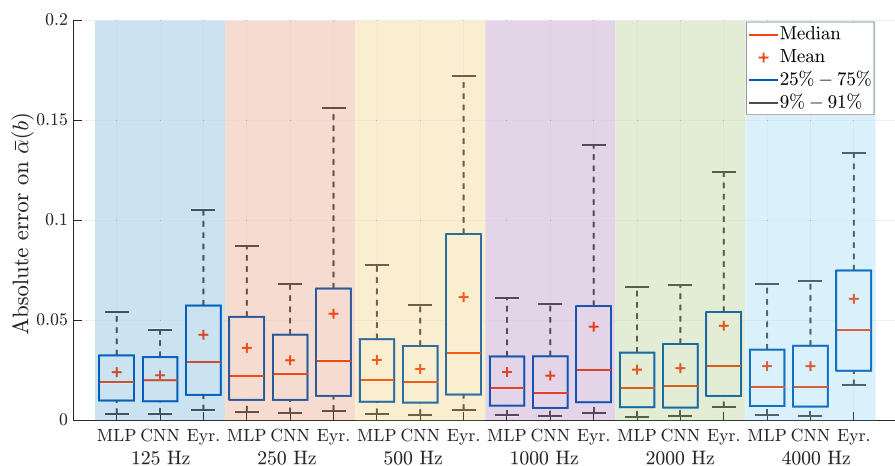


FIG. 7. (Color online) The comparison of $\bar{\alpha}(b)$ estimation errors on the realistic test set in different octave bands. The set used for training networks is RB.

the different models against various acoustical conditions. Unless stated otherwise, the acoustic parameters follow the RB sampling (see Sec. III B), and the RIRs have undergone the same pretreatment as those in Sec. IV A. First, Fig. 6(b) compares results on three test sets, respectively, containing only cube-like rooms ($L_x, L_y \in [2, 4]$; $L_z = 2.5$), flat rooms ($L_x, L_y \in [8, 10]$; $L_z = 2.5$), and elongated rooms ($L_x \in [2, 4]$; $L_y \in [8, 10]$; $L_z = 2.5$). Unsurprisingly, with Eyring’s model, the smallest absolute errors are obtained on cube-like rooms for which the sound field is closest to being diffuse (Hodgson, 1994, 1996). Logically, both the mean and variance of this error increase for the two other geometrical configurations. Although learned models only provide minor improvements over Eyring’s formula under cube-like geometries where the DSF assumption is mostly met, they offer a clear advantage in nonhomogeneous conditions.

Figure 6(c) compares the results for three test sets, each associated with a specific reverberation (slightly reverberant, semi-reverberant, reverberant). Where the obtained errors tend to increase as the reverberation time decreases, the learned models remain superior to Eyring’s in all conditions. For Eyring, this increase is expected as more reverberant rooms are closer to the DSF hypothesis (Hodgson, 1994, 1996).

Figure 6(d) reports errors as a function of the SNR when additive white Gaussian noise is added to the RIR signals (SNR levels are calculated on the first 500 ms of the RIRs). It can be seen that the Eyring model estimations degrade abruptly for SNRs of 30 dB or lower. To investigate this effect, Fig. 8 shows the 1 kHz Schroeder curves of an example RIR under varying noise levels. As can be seen, as the noise level increases, a clean, linear, -30 dB log-energy decay may no longer be available, thus, degrading the RT_{30} estimation. This is a well-known limitation of the reverberation-based techniques, which often require manual adaptation of the decay level used, depending on the measurements. On the other hand, the learned MLP-RB and CNN-RB models, trained on a noisy dataset (30 dB SNR), prove to be much more robust to noise, suggesting that they adaptively extract relevant cues from the RIRs.

Finally, Figs. 6(e) and 6(f) report errors as a function of $\bar{\alpha}$ and the mean scattering coefficient \bar{s} , where each

coefficient is fixed to a constant value across all octave bands and surfaces in each test set. Once again, the behavior of Eyring’s model matches the model expected from the reverberation theory because rooms containing high-scattering, low-absorption materials tend to feature more DSFs (Hodgson, 1994, 1996). On the other hand, learned models perform similarly or better than Eyring’s model for $\bar{s} < 0.5$ and $\bar{\alpha} < 0.5$ but significantly less well otherwise. This is because the mean scattering values outside those ranges were not present in the RB training set [see Fig. 2(c)]. Although learning-based methods show remarkable interpolation capabilities, they are known to have limited extrapolation capabilities.

To get further insight on the influence of the scattering coefficients and diffusion when training neural networks, we tried retraining the CNN model on a purely specular RB set, i.e., using only the image source method in Roomsim while

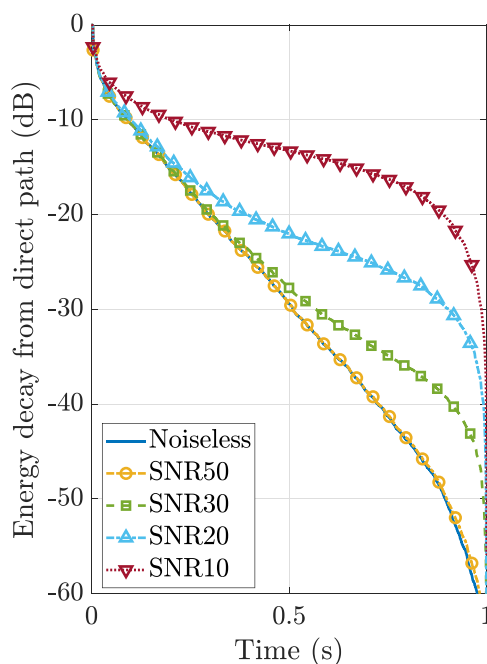


FIG. 8. (Color online) The 1 kHz Schroeder curves of a RIR under varying SNRs is depicted.

disabling the diffuse-rain algorithm as was performed in, e.g., the learning-based absorption estimation technique proposed in Yu and Kleijn (2021). The obtained mean absolute error on $\bar{\alpha}$ of the realistic test set was 0.18, which is six times larger than when using the original RB set with the diffusion activated (0.03). This strongly highlights the importance of taking into account scattering effects when training learning-based acoustic estimation techniques.

Overall, this extensive simulated study reveals that carefully trained, virtually supervised models can consistently and significantly outperform conventional reverberation-based techniques in the task of estimating the quantity $\bar{\alpha}$, particularly under noisy or non-DSF conditions. This was expected as the use of Eyring’s model is theoretically inadequate under such conditions even if the observed absolute errors were reasonable in practice (see Sec. V A). In conditions close to the DSF hypothesis, the learned models and reverberation-based models become comparable. This suggests that the trained models learned a correction with respect to the classical models under non-DSF conditions by extracting richer features from the RIRs than from the mere reverberation times.

VI. TEST ON REAL DATA

A. Real dataset

To evaluate the generalizability of the proposed approach to the real measured RIR, we use a subset of the dEchorate dataset (Di Carlo et al., 2021). The dataset consists of RIR measured in a 6 m × 6 m × 2.4 m acoustic room in the acoustic laboratory of the Bar-Ilan University. The wall and ceiling absorption properties can be changed by flipping the double-sided panels with one reflective and one absorbing face.

Ten different room configurations are considered. They are represented as binary strings of 6 bits in Table I, where “1” denotes a reflective surface, “0” denotes an absorbing surface, and the ordered bits represent the floor, the ceiling, and the West, South, East, and North walls. For each configuration, 90 RIRs from all combinations of 3 sources and 30 receivers spread inside the room are measured. The sources are Avantone Pro Active Mixcube loudspeakers (directional; Middletown, NY), and the receivers are AKG CK32 omnidirectional microphones (Vienna, Austria). Whereas room configurations 1–9 only contain the sources and receivers, room configuration 10 also contains some typical meeting room furnitures, namely, a table, some chairs, and a coat hanger. Each RIR is measured using the exponential sine sweep

technique described in Farina (2007). In this experiment, the octave bands centered at 125 and 250 Hz will not be considered because the measured RIRs did not exhibit sufficient power in those bands for reliable $RT(b)$ estimations. This observation is consistent with the frequency response provided by the loudspeakers’ manufacturer, which decays exponentially from 200 Hz downward.

B. Reference absorption values

A major difficulty in evaluating the considered models on real *in situ* measures is the unavailability of ground truth for the mean absorption coefficients, which would require knowing the true absorption profiles of every material in the room. Although some of them could be inferred from the manufacturer’s data, only coarse values of $\bar{\alpha}(b)$ would be obtained in this way. To overcome this difficulty while ensuring that a single, stable, and reliable mean absorption profile is used as a reference for each room, we propose a technique based on the aggregation of multiple RIR measurements.

For each room configuration, the Schroeder curves of the 90 measured RIRs in 4 octave bands were traced (Schroeder, 1965). Then, the Schroeder curves were visually inspected and separated into two sets. Set \mathcal{A} contains Schroeder curves featuring sufficient linear log-energy decay from -5 to -15 dB at least. Set \mathcal{B} contains all of the other curves. In practice, 49% of the 3600 Schroeder curves were discarded to the set \mathcal{B} in this way. These mostly corresponded to challenging measurement situations contained in the dEchorate dataset, such as a receiver near a surface or a loudspeaker facing toward a surface and away from receivers. Then, for each room configuration and each octave band b , the reference mean absorption coefficient $\bar{\alpha}_{\text{ref}}(b)$ is taken to be the median value of Eyring’s model based on the $RT_{10}(b)$ computed from the Schroeder curves in \mathcal{A} only and the known room’s volume and total surface. This median value $\bar{\alpha}_{\text{ref}}(b)$ is taken over at least 5 and, on average, 47 estimates (see Table I), yielding a reliable and robust value. As can be seen in Table I, a diversity of the mean absorption coefficients $\bar{\alpha}_{\text{ref}}(b)$ between 0.12 and 0.52 is represented. This matches quite well with the range of values considered in this study [see Figs. 1 and 2(c)].

To further validate this choice of reference value, the left part of Fig. 9 shows the means and standard deviations (stds) of the absolute differences between the single-RIR Eyring estimates and the proposed median-based reference for each room configuration and each octave band using

TABLE I. Absorption coefficients $\bar{\alpha}_{\text{ref}}(b)$ calculated in the ten room configurations. For each coefficient, the number of corresponding Schroeder curves in \mathcal{A} used to compute the median Eyring’s estimate is given in parentheses. Room 10 contains furniture.

	Room 1	Room 2	Room 3	Room 4	Room 5	Room 6	Room 7	Room 8	Room 9	Room 10
Configuration	000000	011000	011100	011110	011111	001000	000100	000010	000001	010001
500 Hz	0.42 (11)	0.23 (7)	0.20 (20)	0.17 (51)	0.13 (48)	0.39 (8)	0.38 (5)	0.40 (8)	0.35 (7)	0.23 (12)
1000 Hz	0.52 (62)	0.28 (83)	0.25 (86)	0.17 (89)	0.13 (90)	0.44 (79)	0.41 (74)	0.44 (69)	0.43 (70)	0.33 (72)
2000 Hz	0.50 (65)	0.34 (81)	0.30 (86)	0.19 (82)	0.14 (88)	0.44 (74)	0.42 (64)	0.44 (66)	0.44 (67)	0.37 (69)
4000 Hz	0.37 (15)	0.35 (17)	0.29 (22)	0.16 (16)	0.12 (29)	0.38 (17)	0.33 (12)	0.32 (14)	0.34 (18)	0.32 (14)

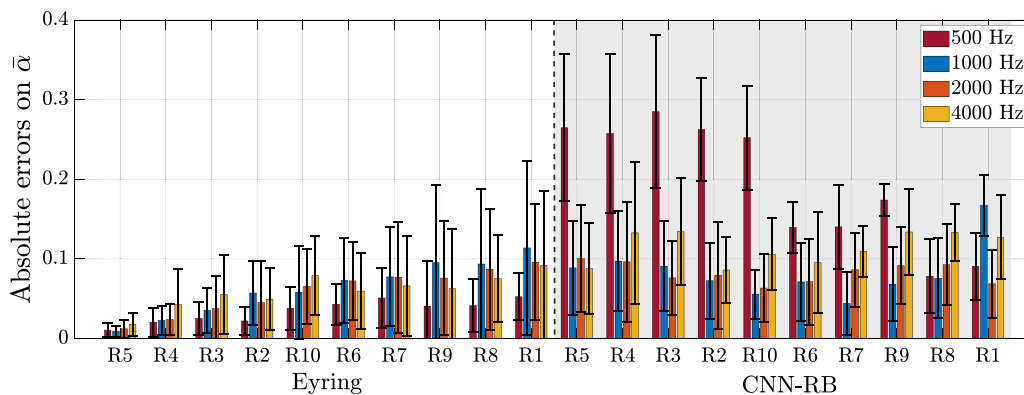


FIG. 9. (Color online) The comparison of the $\bar{\alpha}(b)$ mean estimation errors over measured RIRs in 10 rooms and 4 octave bands with Eyring and CNN-RB. Only selected RIRs with Schroeder curves in \mathcal{A} are included.

RIRs from set \mathcal{A} only. The rooms are sorted from left-to-right from the most reverberant room to the least reverberant room. Clearly, it appears that both the means and stds of the differences between the single and median-based estimates increase as the reverberation time decreases consistently with the reverberation theory (Hodgson, 1994, 1996). Nevertheless, both these means and stds remain reasonably low (below 0.1) under all configurations, despite measurements being taken from many different source-receiver placements in the room. This validates our premise of a close-to-DSF in these experiments, at least when restricting to the RIRs inside of the set \mathcal{A} for each octave band.

C. Real data results

On real RIRs, the MLP models appeared to perform significantly worse than the CNN models, yielding errors up to twice as large. This is consistent with the better generalization capabilities of the CNN models observed in Fig. 5 and discussed in Sec. IV B. We, hence, omit the MLP results in the remainder of this section for compactness.

The right part of Fig. 9 reports the mean and stds of the absolute errors for the CNN-RB model using only the RIRs in \mathcal{A} . Encouragingly, for the 1, 2, and 4 kHz octave bands, the learning-based method yields errors below or around 0.1 for all rooms, which is a reasonable uncertainty in the context of the acoustic diagnosis. The errors are comparable to those obtained with Eyring’s formula except in the three most reverberant rooms ($R3, R4,$ and $R5$) for which the latter performs very well. For the octave band centered at 4 kHz, the CNN-RB errors increase slightly. A possible explanation could lie in the stronger directivity of the source at this frequency as observed in the manufacturer’s data (recall that the neural network has only been trained on omnidirectional sources). For the octave band centered at 500 Hz, the CNN-RB errors are much larger in all of the rooms except $R1$ and $R8$. One of the preferred hypotheses is the existence of a wave phenomenon in this band that could not be learned by the neural network trained on Roomsim. These hypotheses will need to be validated by further research on real data. Figure 10 shows the same results in the form of bar plots for the 1 kHz octave band, further confirming that the CNN-RB

model yields error distributions comparable to those of Eyring in this band.

Finally, Fig. 11 compares errors obtained with the CNN-RB on measured RIRs whose 1 kHz Schroeder curves are in \mathcal{A} against those whose Schroeder curves are in \mathcal{B} . Note that rooms $R3, R4,$ and $R5$ are omitted here because an insufficient number of curves were placed in \mathcal{B} for these rooms. Encouragingly, we observe that the CNN is largely unaffected by the nonlinear or insufficient log-energy decays of Schroeder curves in \mathcal{B} . This suggests that the network learned to rely on more elaborate and more robust features than those used by the reverberation-based techniques. In contrast, obtaining reliable absorption estimates from these curves using Eyring’s model was fundamentally impossible due to its reliance on the reverberation time.

VII. CONCLUSION

In this work, we tackled the inverse problem of estimating the area-weighted mean absorption coefficients of a room from a single RIR using virtually supervised learning in a broad range of acoustical conditions pertaining to the field of building acoustic diagnoses. Different neural network designs and simulated training strategies were proposed, explored, and tested. The developed methods were

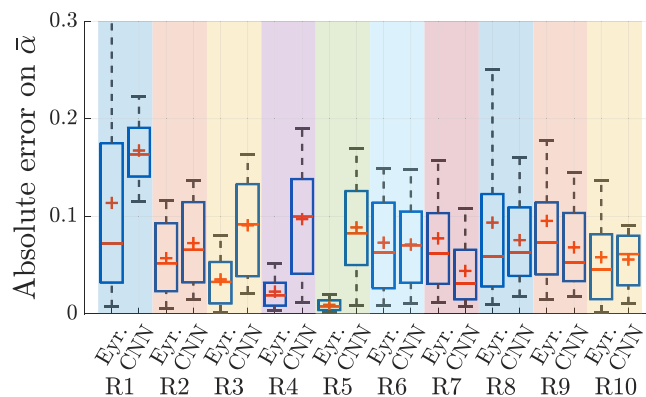


FIG. 10. (Color online) The comparison of the $\bar{\alpha}(1000 \text{ Hz})$ estimation errors over measured RIRs in ten rooms using Eyring and CNN-RB. Only selected RIRs with Schroeder curves in \mathcal{A} are included.

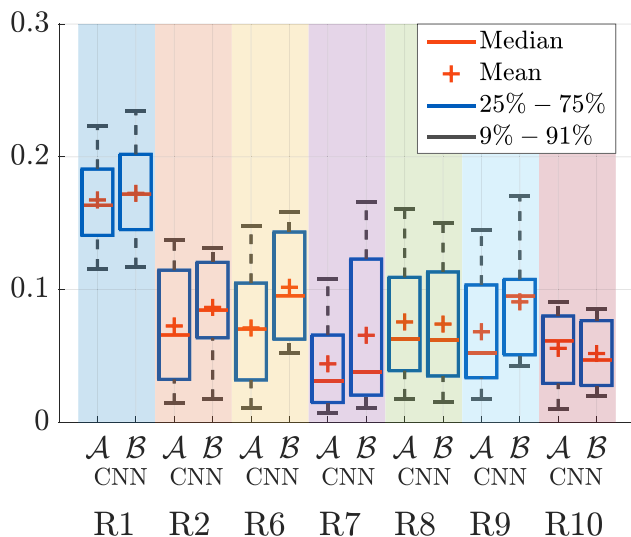


FIG. 11. (Color online) The comparison of the $\bar{\alpha}(1000\text{ Hz})$ CNN-RB estimation errors over measured RIRs with 1 kHz Schroeder curves in \mathcal{A} vs those in \mathcal{B} .

compared to classical formulas that hinge on the room’s volume, total surface, reverberation time, and DSF hypothesis. In close-to-DSF conditions, our experiments on both simulated and real data revealed that the best learned models yielded estimation errors comparable to classical errors without needing the room’s geometry. As expected and predicted by the reverberation theory, the performances of the DSF-based models degraded under conditions departing from the DSF. These include rooms featuring less reverberation, less diffusion, non-homogenous geometries, and, more generally, RIRs featuring insufficient or nonlinear decays of their Schroeder curves. In contrast, the proposed virtually trained models showed remarkable robustness in estimating the target quantity under such conditions, suggesting that they learned to rely on more elaborate and more robust features than those used by the reverberation-based techniques.

This first extensive experimental study on virtually supervised mean absorption estimation aimed at paving the way toward simpler and more robust acoustic diagnosis techniques. Future work will include further experimental investigations on the poorer performance of the learned models at lower frequencies on real data, notably by employing higher-end sound sources. Leads for improving the learned models include domain adaption, data augmentation, and probabilistic uncertainty modeling. We also plan to build on our findings to tackle the much more difficult problem of estimating the absorption coefficients of individual surfaces from the RIRs. For this, geometrically informed models and the aggregation of the RIRs from multiple source-receiver pairs will be leveraged.

¹The full lists of materials and associated absorption profiles considered in this study are available at https://members.loria.fr/ADeleforge/files/jasa2021_supplementary_material.zip (Last viewed August 16, 2021).

²We denote by $RT_x(b)$ a reverberation time calculated on a Schroeder curve’s slope from -5 to $-5 - X$ dB (Schroeder, 1965).

Allard, J. F., and Akinine, A. (1985). “Acoustic impedance measurements with a sound intensity meter,” *Appl. Acoust.* **18**, 69–75.

Allard, J. F., and Sieben, B. (1985). “Measurements of acoustic impedance in a free field with two microphones and a spectrum analyzer,” *J. Acoust. Soc. Am.* **77**, 1617–1618.

Allen, J. B., and Berkley, D. A. (1979). “Image method for efficiently simulating small-room acoustics,” *J. Acoust. Soc. Am.* **65**(4), 943–950.

Ando, Y. (1968). “The interference pattern method of measuring the complex reflection coefficient of acoustic materials at oblique incidence,” in *Proc. 6th International Congress on Acoustics (ICA, Tokyo, Japan)*.

Aoshima, N. (1981). “Computer-generated pulse signal applied for sound measurement,” *J. Acoust. Soc. Am.* **69**, 1484–1488.

ASTM (2006). E1050-98. “Standard test method for impedance and absorption of acoustical materials using a tube, two microphones, and a digital frequency analysis system” (American Society for Testing and Materials, Philadelphia, PA).

Barry, T. (1974). “Measurement of the absorption spectrum using correlation/spectral density techniques,” *J. Acoust. Soc. Am.* **55**, 1349–1351.

Bianco, M. J., Gerstoft, P., Traer, J., Ozanich, E., Roch, M. A., Gannot, S., and Deledalle, C.-A. (2019). “Machine learning in acoustics: Theory and applications,” *J. Acoust. Soc. Am.* **146**(5), 3590–3628.

Borish, J. (1984). “Extension of the image model to arbitrary polyhedra,” *J. Acoust. Soc. Am.* **75**(6), 1827–1836.

Botteldooren, D. (1995). “Finite-difference time-domain simulation of low-frequency room acoustic problems,” *J. Acoust. Soc. Am.* **98**(6), 3302–3308.

Brandão, E., Lenzi, A., and Paul, S. (2015). “A review of the *in situ* impedance and sound absorption measurement techniques,” *Acta Acust. Acust.* **101**(3), 443–463.

Chakrabarty, S., and Habets, E. A. (2017). “Broadband doa estimation using convolutional neural networks trained with noise signals,” in *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)* (IEEE, New York), pp. 136–140.

Champoux, Y., and L’espérance, A. (1988). “Numerical evaluation of errors associated with the measurement of acoustic impedance in a free field using two microphones and a spectrum analyzer,” *J. Acoust. Soc. Am.* **84**, 30–38.

Champoux, Y., Nicolas, J., and Allard, J. F. (1988). “Measurement of acoustic impedance in a free field at low frequencies,” *J. Sound Vib.* **125**, 313–323.

Chung, J., and Blaser, D. (1980a). “Transfer function method of measuring in-duct acoustic properties. I. Theory,” *J. Acoust. Soc. Am.* **68**(3), 907–913.

Chung, J., and Blaser, D. (1980b). “Transfer function method of measuring in-duct acoustic properties. II: Experiment,” *J. Acoust. Soc. Am.* **68**(3), 914–921.

Cramond, A. J., and Don, C. G. (1984). “Reflection of impulses as a method of determining acoustic impedance,” *J. Acoust. Soc. Am.* **75**, 382–389.

Davies, J. C., and Mulholland, K. A. (1979). “An impulse method of measuring normal impedance at oblique incidence,” *J. Sound Vib.* **67**, 135–149.

Deecke, V. B., and Janik, V. M. (2006). “Automated categorization of bioacoustic signals: Avoiding perceptual pitfalls,” *J. Acoust. Soc. Am.* **119**(1), 645–653.

Deleforge, A., Forbes, F., and Horaud, R. (2015a). “Acoustic space learning for sound-source separation and localization on binaural manifolds,” *Int. J. Neural Syst.* **25**(01), 1440003.

Deleforge, A., Horaud, R., Schechner, Y. Y., and Girin, L. (2015b). “Co-localization of audio sources in images using binaural features and locally-linear regression,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.* **23**(4), 718–731.

Di Carlo, D., Deleforge, A., and Bertin, N. (2019). “MIRAGE: 2D source localization using microphone pair augmentation with echoes,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE, New York), pp. 775–779.

Di Carlo, D., Tandeitnik, P., Foy, C., Deleforge, A., Bertin, N., and Gannot, S. (2021). “dechorate: A calibrated room impulse response database for echo-aware signal processing,” *arXiv:2104.13168*.

Farina, A. (2000). “Simultaneous measurement of impulse response and distortion with a swept-sine technique,” in *Audio Engineering Society Convention 108* (Audio Engineering Society, New York).

- Farina, A. (2007). "Advancements in impulse response measurements by sine sweeps," in *Audio Engineering Society Convention 122* (Audio Engineering Society, New York).
- Gamper, H., and Tashev, I. J. (2018). "Blind reverberation time estimation using a convolutional neural network," in *2018 16th International Workshop on Acoustic Signal Enhancement (IWAENC)* (IEEE, New York), pp. 136–140.
- Garai, M. (1993). "Measurement of the sound-absorption coefficient in situ: The reflection method using periodic pseudorandom sequences of maximum length," *Appl. Acoust.* **39**, 119–139.
- Gaultier, C., Kataria, S., and Deleforge, A. (2017). "Vast: The virtual acoustic space traveler dataset," in *International Conference on Latent Variable Analysis and Signal Separation* (Springer, New York), pp. 68–79.
- Genovese, A. F., Gamper, H., Pulkki, V., Raghuvanshi, N., and Tashev, I. J. (2019). "Blind room volume estimation from single-channel noisy speech," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE, New York), pp. 231–235.
- Gradišek, A., Slapničar, G., Šorn, J., Luštrek, M., Gams, M., and Grad, J. (2017). "Predicting species identity of bumblebees through analysis of flight buzzing sounds," *Bioacoustics* **26**(1), 63–76.
- Guidorzia, P., Barbaresia, L., D'Orazio, D., and Garai, M. (2015). "Impulse responses measured with MLS or swept-sine signals applied to architectural acoustics: An in-depth analysis of the two methods and some case studies of measurements inside theaters," in *6th International Building Physics Conference* (Elsevier, Amsterdam, Netherlands).
- Habets, E. A. (2006). "Room impulse response generator," Technische Universiteit Eindhoven, Tech. Rep. 2(2.4), 1.
- He, W., Motlicek, P., and Odobez, J.-M. (2019). "Adaptation of multiple sound source localization neural networks with weak supervision and domain-adversarial training," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE, New York), pp. 770–774.
- Hodgson, M. (1994). "When is diffuse-field theory accurate?," *Can. Acoust.* **22**(3), 41–42.
- Hodgson, M. (1996). "When is diffuse-field theory applicable?," *Appl. Acoust.* **49**(3), 197–201.
- Hollin, K. A., and Jones, M. H. (1977). "The measurement of sound absorption coefficient *in situ* by a correlation technique," *Acustica* **37**, 103–110.
- Ingård, U., and Bolt, R. H. (1951). "A free field method of measuring the absorption coefficient of acoustic materials," *J. Acoust. Soc. Am.* **23**, 509–516.
- ISO (2001). 10534:2001. "Acoustics. Determination of sound absorption coefficient and impedance in impedance tubes. Part 1: Method using standing wave. Part 2: Transfer function method" (International Organization for Standardization, Geneva, Switzerland).
- ISO (2003). 354:2003. "Acoustics—Measurement of sound absorption in a reverberation room" (International Organization for Standardization, Geneva, Switzerland).
- ISO (2008). 3382-2:2008. "Acoustics—Measurement of room acoustic parameters—Part 2: Reverberation time in ordinary rooms" (International Organization for Standardization, Geneva, Switzerland).
- Kataria, S., Gaultier, C., and Deleforge, A. (2017). "Hearing in a shoe-box: Binaural source position and wall absorption estimation using virtually supervised learning," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE, New York), pp. 226–230.
- Kim, C., Misra, A., Chin, K., Hughes, T., Narayanan, A., Sainath, T., and Bacchiani, M. (2017). "Generation of large-scale simulated utterances in virtual rooms to train deep-neural networks for far-field speech recognition in google home," *Interspeech 2017* (ISCA, Stockholm, Sweden), pp. 379–383.
- Kingma, D. P., and Ba, J. (2014). "Adam: A method for stochastic optimization," [arXiv:1412.6980](https://arxiv.org/abs/1412.6980).
- Kulowski, A. (1985). "Algorithmic representation of the ray tracing technique," *Appl. Acoust.* **18**(6), 449–469.
- Kuttruff, H. (2009). *Room Acoustics*, 5th ed. (Spon, Oxfordshire, England).
- Lefort, R., Real, G., and Drémeau, A. (2017). "Direct regressions for underwater acoustic source localization in fluctuating oceans," *Appl. Acoust.* **116**, 303–310.
- Li, J. F., and Hodgson, M. (1997). "Use of pseudo-random sequences and a single microphone to measure surface impedance at oblique incidence," *J. Acoust. Soc. Am.* **102**, 2200–2210.
- Luo, Y., and Mesgarani, N. (2018). "Tasnet: Time-domain audio separation network for real-time, single-channel speech separation," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE, New York), pp. 696–700.
- Mesaros, A., Heittola, T., Diment, A., Elizalde, B., Shah, A., Vincent, E., Raj, B., and Virtanen, T. (2017). "Dcase 2017 challenge setup: Tasks, datasets and baseline system," in *DCASE 2017-Workshop on Detection and Classification of Acoustic Scenes and Events*.
- Mesaros, A., Heittola, T., and Virtanen, T. (2019). "Acoustic scene classification in dcase 2019 challenge: Closed and open set classification and data mismatch setups,".
- Mintén, M., Cops, A., and Lauriks, W. (1988). "Absorption characteristics of an acoustic material at oblique incidence measured with the two-microphone technique," *J. Sound Vib.* **120**, 499–510.
- Müller, S., and Massarani, P. (2001). "Transfer-function measurement with sweeps. Director's cut including previously unreleased material and some corrections," *J. Audio Eng. Soc.* **49**(6), 443–471.
- Niu, H., Reeves, E., and Gerstoft, P. (2017). "Source localization in an ocean waveguide using supervised machine learning," *J. Acoust. Soc. Am.* **142**(3), 1176–1188.
- Nobile, M. A., and Hayek, S. I. (1985). "Acoustic propagation over an impedance plane," *J. Acoust. Soc. Am.* **78**, 1325–1336.
- Nolan, M. (2020). "Estimation of angle-dependent absorption coefficients from spatially distributed *in situ* measurements," *J. Acoust. Soc. Am.* **147**(2), EL119–EL124.
- Nolan, M., Fernandez-Grande, E., Brunskog, J., and Jeong, C.-H. (2018). "A wavenumber approach to quantifying the isotropy of the sound field in reverberant spaces," *J. Acoust. Soc. Am.* **143**(4), 2514–2526.
- Okuzono, T., Otsuru, T., Tomiku, R., and Okamoto, N. (2014). "A finite-element method using dispersion reduced spline elements for room acoustics simulation," *Appl. Acoust.* **79**, 1–8.
- Parsons, S., and Jones, G. (2000). "Acoustic identification of twelve species of echolocating bat by discriminant function analysis and artificial neural networks," *J. Exp. Biol.* **203**(17), 2641–2656.
- Peterson, P. M. (1986). "Simulating the response of multiple microphones to a single acoustic source in a reverberant room," *J. Acoust. Soc. Am.* **80**(5), 1527–1529.
- Pietrzyk, A. (1998). "Computer modeling of the sound field in small rooms," in *Audio Engineering Society Conference: 15th International Conference: Audio, Acoustics & Small Spaces* (Audio Engineering Society, New York).
- Prawda, K., Schlecht, S. J., and Välimäki, V. (2020). "Evaluation of reverberation time models with variable acoustics," in *18th Sound and Music Computing Conference* (Axa sas/SMC Network, Torino, Italy).
- Rathsam, J., and Rafaely, B. (2015). "Analysis of absorption *in situ* with a spherical microphone array," *Appl. Acoust.* **89**, 273–280.
- Richard, A., Fernandez-Grande, E., Brunskog, J., and Jeong, C.-H. (2017). "Estimation of surface impedance at oblique incidence based on sparse array processing," *J. Acoust. Soc. Am.* **141**(6), 4115–4125.
- Rife, D., and Vanderkooy, J. (1999). "Transfer-function measurement with maximum length sequences," *J. Audio Eng. Soc.* **37**(6), 419–444.
- Samarasinghe, P. N., Abhayapala, T. D., Lu, Y., Chen, H., and Dickins, G. (2018). "Spherical harmonics based generalized image source method for simulating room acoustics," *J. Acoust. Soc. Am.* **144**(3), 1381–1391.
- Scheibler, R., Bezzam, E., and Dokmanić, I. (2018). "Pyroomacoustics: A python package for audio room simulation and array processing algorithms," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE, New York), pp. 351–355.
- Schimmel, S. M., Muller, M. F., and Dillier, N. (2009). "A fast and accurate 'shoebox' room acoustics simulator," in *2009 IEEE International Conference on Acoustics, Speech and Signal Processing* (IEEE, New York), pp. 241–244.
- Schröder, D. (2011). *Physically Based Real-Time Auralization of Interactive Virtual Environments* (Logos, Berlin), Vol. 11.
- Schroeder, M. R. (1965). "New method of measuring reverberation time," *J. Acoust. Soc. Am.* **37**, 409–412.
- Schroeder, M. R. (1979). "Integrated-impulse method measuring sound decay without using impulses," *J. Acoust. Soc. Am.* **66**(2), 497–500.
- Schroeder, M. R. (1996). "The Schroeder frequency revisited," *J. Acoust. Soc. Am.* **99**, 3240–3241.
- Sides, D. J., and Mulholland, K. A. (1971). "The variation of normal layer impedance with angle of incidence," *J. Sound Vib.* **14**, 139–142.

- Stan, G.-B., Embrechts, J.-J., and Archambeau, D. (2002). "Comparison of different impulse response measurement techniques," *J. Audio Eng. Soc.* **50**(4), 249–262.
- Suzuki, Y., Asano, F., Kim, H., and Sone, T. (1995). "An optimum computer-generated pulse signal suitable for the measurement of very long impulse responses," *J. Acoust. Soc. Am.* **97**(2), 1119–1123.
- Tamura, M. (1990). "Spatial Fourier transform method of measuring reflection coefficients at oblique incidence. I: Theory and numerical examples," *J. Acoust. Soc. Am.* **88**(5), 2259–2264.
- Torras-Rosell, A., and Jacobsen, F. (2010). "Measuring long impulse responses with pseudorandom sequences and sweep signals," in *Internoise*, Lisbon, Portugal.
- Vorländer, M., and Mommertz, E. (2000). "Definition and measurement of random-incidence scattering coefficients," *Appl. Acoust.* **60**(2), 187–199.
- Wabnitz, A., Epain, N., Jin, C., and Van Schaik, A. (2010). "Room acoustics simulation for multichannel microphone arrays," in *Proceedings of the International Symposium on Room Acoustics* (ICA, Melbourne, Australia), pp. 1–6.
- Yu, W., and Kleijn, W. B. (2021). "Room acoustical parameter estimation from room impulse responses using deep neural networks," *IEEE/ACM Trans. Audio, Speech, Lang. Process.* **29**, 436–447.
- Yuzawa, M. (1975). "A method of obtaining the oblique incident sound absorption coefficient through an on-the-spot measurement," *Appl. Acoust.* **8**, 27–41.