



**HAL**  
open science

# A Robust Denoising Process for Directional Room Impulse Responses with Diffuse Reverberation Tails

Pierre Massé, Thibaut Carpentier, Olivier Warusfel, Markus Noisternig

► **To cite this version:**

Pierre Massé, Thibaut Carpentier, Olivier Warusfel, Markus Noisternig. A Robust Denoising Process for Directional Room Impulse Responses with Diffuse Reverberation Tails. 2020. hal-02443679v1

**HAL Id: hal-02443679**

**<https://hal.science/hal-02443679v1>**

Preprint submitted on 17 Jan 2020 (v1), last revised 20 Apr 2020 (v2)

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

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

# **A Robust Denoising Process for Directional Room Impulse Responses with Diffuse Reverberation Tails**

Pierre Massé,<sup>1</sup> Thibaut Carpentier,<sup>1</sup> Olivier Warusfel,<sup>1</sup> and Markus Noisternig<sup>1</sup>

*Acoustic and Cognitive Spaces group, Sorbonne Université, IRCAM, CNRS, STMS,  
75004 Paris, France<sup>a)</sup>*

1 Directional room impulse responses (DRIR) measured using spherical microphone  
2 arrays (SMA) are seeing increasingly widespread use in reproducing room reverber-  
3 ation effects on three-dimensional surround sound systems (e.g. Higher-Order Am-  
4 bisonics) through multi-channel DRIR convolution. However, such measured im-  
5 pulse responses inevitably present a non-negligible noise floor, which may lead to  
6 a perceptible “infinite reverberation effect” when convolved with an input sound.  
7 Furthermore, individual sensor noise and momentary measurement artefacts may ad-  
8 ditionally corrupt the resulting impulse response. This paper presents a robust DRIR  
9 denoising procedure applicable to impulse responses with diffuse late reverberation  
10 tails, which can be modeled by a stochastic process. In such cases, the non-decaying  
11 frequency-dependent noise floor may be replaced by a synthesized diffuse tail param-  
12 eterized by the DRIR’s energy decay envelope. It is shown that performing such tail  
13 re-synthesis in the spherical harmonic domain (SHD), using an independent zero-  
14 mean Gaussian noise for each component, preserves not only the reverberation tail’s  
15 frequency-dependent decay properties, but also its spatial incoherence. The proposed  
16 process is then evaluated through its application to DRIRs measured in real-world  
17 conditions, and finally some aspects of performance and consistency verification are  
18 discussed.

---

<sup>a)</sup> [pierre.masse@ircam.fr](mailto:pierre.masse@ircam.fr);

19 **I. INTRODUCTION**

20 **A. SMA analysis in the SHD**

21 Spherical microphone arrays (SMA) enable the directional analysis of a given sound  
 22 field by sampling it over the  $Q$  transducer positions on their surface. A natural choice  
 23 of representation for a function defined on such a surface  $S^2$  is the spherical harmonic  
 24 domain (SHD), whose basis functions  $Y_{l,m}$  are analogues of the trigonometric functions in  
 25 the application of Fourier expansion theory on the sphere ([Driscoll and Healy, 1994](#)):

$$X_{l,m}(f, t) = \int_{\Omega \in S^2} x(f, \Omega, t) Y_{l,m}(\Omega) d\Omega, \tag{1}$$

26 where  $\Omega = (\theta, \phi)$  is a point on the surface of a sphere with fixed radius  $r = a$  (in  
 27 conformity with [ISO8000-2:2009 \(E\)](#)),  $x(f, \Omega, t)$  is the time-frequency domain representation  
 28 of the sound field on the sphere, and  $Y_{l,m}(\Omega)$  are the spherical harmonics of order  $l \in \mathbb{Z}^+$   
 29 and degree  $m \in [-l, l]$ . This transform thus defines the SHD signal coefficients  $X_{l,m}(f, t)$  for  
 30 each component or mode  $(l, m)$ . Using an SMA, the integral in Eq. (1) is discretized and  
 31 can be approximated by a weighted sum over the microphone positions; the specific weights  
 32 are chosen such that the sum approaches the ideal integral of Eq. (1), e.g. by least-squares  
 33 minimization ([Rafaely, 2005](#)).

34 The discrete transform can be simply written in matrix form:

$$\mathbf{x}_{\text{SHD}}(f, t) = \mathbf{Y}\mathbf{x}(f, t), \tag{2}$$

35 where  $\mathbf{x}(f, t)$  is the column vector containing the time-frequency representation of the  
 36 signal measured at each transducer position  $\Omega_q$ ,  $\mathbf{Y}$  is the  $(L + 1)^2 \times Q$  encoding matrix of  
 37 elements  $y_{q,n} = \alpha_q Y_{l,m}(\Omega_q)$  (with indices  $n = l^2 + l + m + 1$  up to a maximum encoding  
 38 order  $L$ , and  $\alpha_q$  the aforementioned array weights), and  $\mathbf{x}_{\text{SHD}}(f, t)$  is the column vector of  
 39 resulting SHD coefficients. The array’s sampling configuration must then lead to an encod-  
 40 ing matrix with  $K$  non-vanishing singular values such that  $K = (L + 1)^2 \leq Q$  (Noisternig  
 41 *et al.*, 2011), thereby effectively limiting the maximum achievable order  $L$  for a given SMA.  
 42 Finally, in order to obtain an array-independent representation of the measured sound field,  
 43 a subsequent correction for the so-called mode strengths (or holographic functions) of the  
 44 SMA must be applied. Such is the case in the widespread Higher-Order Ambisonics (HOA)  
 45 format, where the center of the sphere is used as the reference point and for which the  
 46 correcting filters are determined accordingly (Daniel and Moreau, 2004).

## 47 B. Previous work

48 Monophonic room impulse responses (RIR) have long been modelled as an exponentially  
 49 decaying stochastic process (Schroeder, 1962), which has been shown to be valid assuming  
 50 sufficiently high echo density and modal overlap is achieved (Polack, 1988). These conditions  
 51 lead to a lower time limit for echo density, known as the “mixing time”, and a lower frequency  
 52 limit for modal overlap, known as the “Schroeder frequency”. Beyond these limits, the late  
 53 reverberation field is considered to be fully “diffuse”, i.e. it behaves as a spatially isotropic  
 54 distribution of a statistically significant number of incoherent and uncorrelated plane waves.  
 55 Such a field can be synthesized in the form of a zero-mean Gaussian noise filtered by an

56 exponentially-decaying energy envelope (Jot *et al.*, 1997). This envelope is parameterized by  
 57 a frequency-dependent decay coefficient  $\delta(f)$  (usually represented as the 60 dB reverberation  
 58 time,  $T_{60}(f) = 3 \ln(10)/\delta(f)$ ) and an initial power spectrum  $P_0(f)$ ; these parameters can  
 59 be extracted by analysis of the energy decay relief (EDR), a time-frequency extension of  
 60 the Schroeder energy decay curve (EDC) (Jot *et al.*, 1997). Non-decaying background noise  
 61 present in a measured impulse response can therefore be replaced by a synthesized zero-mean  
 62 Gaussian noise filtered by a prolongation of the energy decay envelope (Jot *et al.*, 1997). As  
 63 a result, the final signal-to-noise ratio (SNR) is limited only by the quantization noise floor  
 64 for the chosen synthesis bit depth,  $P_{QN} = 20 \log_{10}(2^{-d})$  dB, where  $d$  is the signal bit depth.

65 Guski and Vorländer (2014) have since presented a variety of other noise compensation  
 66 methods, but these focus more on regularizing the broadband EDC calculation in order to  
 67 improve the accuracy of extracted room acoustics parameters (e.g.  $T_{60}$ ,  $C_{80}$  clarity, etc.),  
 68 rather than faithfully re-synthesizing the reverberation tail for convolution applications.  
 69 Some of their techniques resemble that proposed by Cabrera *et al.* (2011) for auralizing  
 70 measured RIRs; all have so far only been presented in the monophonic single-slope decay  
 71 case. Furthermore, eschewing tail re-synthesis for simple decay envelope adjustment places  
 72 strict conditions on the content of the background noise (as noted by Cabrera *et al.* (2011)),  
 73 which may easily not be verified in many “real-world” measurement conditions.

74 Preliminary extensions of Jot’s tail re-synthesis process to the spatialized DRIR case  
 75 were presented by Carpentier *et al.* (2013) (using a reference diffuse field simulated by large  
 76 numbers of incident plane waves to denoise the individual SMA transducer signals) and  
 77 Noisternig *et al.* (2014) (in the SHD using independent zero-mean Gaussian noise realizations

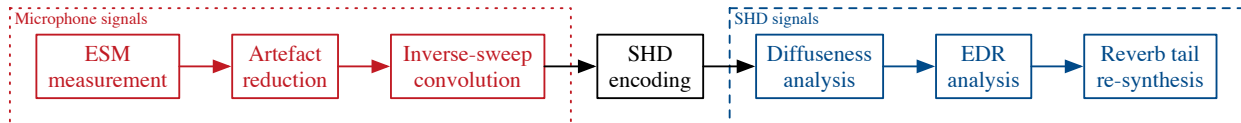


FIG. 1. Outline of the proposed DRIR denoising process, from the initial exponential sweep method (ESM) measurement through to the energy decay relief (EDR) analysis and reverberation tail re-synthesis in the spherical harmonic domain (SHD).

78 per component), both once again in the single-slope decay case. The current work builds  
 79 upon and further details these methods, allowing for multiple-slope decays (such as those  
 80 observed in certain coupled-volume configurations) and demonstrating that tail re-synthesis  
 81 in the SHD guarantees preservation of the late reverberation’s incoherence properties.

## 82 II. PROPOSED DENOISING PROCESS

83 The different parts of the proposed denoising framework are presented in this section, and  
 84 the sequencing of the individual steps is outlined schematically in Fig. 1. The exponential  
 85 sweep method (ESM) measurement and subsequent inverse-sweep convolution are based on  
 86 Farina (2000) and performed on each SMA transducer signal independently; between these  
 87 two steps we introduce an artefact reduction procedure described in section II A. The SHD  
 88 encoding is based on the theory presented in the introduction, and the EDR analysis is  
 89 an extension of Jot *et al.* (1997), detailed in section II B 1. Finally, the main focus of this  
 90 work is on the diffuseness and mixing time analysis (section II B 2) and reverberation tail  
 91 re-synthesis (section II B 3).

## 92    **A. Measurement artefact reduction**

93    Measuring impulse responses using the ESM in so-called “real-world” conditions is in-  
94    evitably subject to three main risk factors: the presence of constant, stationary background  
95    noise (including transducer self-noise), any non-stationarity of the measurement conditions  
96    (temperature, humidity, etc.), and the occurrence of non-stationary noise events. The first is  
97    what is assumed in previous work on the subject and what is aimed to be removed in the  
98    tail re-synthesis procedure. The second can lead to time-variance in the impulse responses  
99    which would require post-processing correction techniques using *a priori* information on the  
100    measurement conditions, and will not be considered in this study. The third is what we will  
101    refer to here as “measurement artefacts”, i.e. short-term sonic events occurring during the  
102    measurement, and reducing their impact is the aim of this section.

103    As noted by [Farina \(2000\)](#), averaging a repetition of several sweeps is a simple way to  
104    increase the SNR, since the ensemble mean of any incoherent stationary noise will tend  
105    to zero as the number of repetitions increases. However, any non-stationary noise events  
106    present in the repetitions will inevitably end up in the noise floor of the average. This is  
107    especially troublesome when considering Schroeder-type reverse-integrated analysis such as  
108    the EDR, since these artefacts will not only accumulate in the reverse-integration of the noise  
109    floor, they will also deviate substantially from the theoretical profile of a reverse-integrated  
110    constant-power noise floor (see section [II B 1](#) below).

111    In an attempt to minimize the influence of these non-stationary noise events, the mag-  
112    nitude spectrograms of the individual sweep repetitions are compared amongst each other



113 in order to identify artefacts, using non-negligible positive deviations from the mean magni-  
 114 tude spectrogram as a discriminating criterion. This maximum allowed deviation is defined  
 115 as  $\xi(f, t) = \mu(f, t) + \alpha\sigma(f, t)$ , where  $\mu(f, t)$  is the mean magnitude spectrogram,  $\sigma(f, t)$  is  
 116 the standard deviation over the available repetitions, and  $\alpha$  is an empirically-set deviation  
 117 factor used as a control parameter. Artefact magnitude values identified as greater than  
 118  $\xi(f, t)$  in each realisation are then replaced with the corresponding mean magnitude over  
 119 the remaining repetitions.

120 This process is applied independently to the ESM measurement signals recorded by each  
 121 SMA transducer. Some example results are illustrated and discussed in section III A.

## 122 B. Reverberation tail analysis and re-synthesis

123 In this section, we first review the energy decay relief (EDR) analysis procedure used  
 124 to extract the reverberation decay parameters, and then present a characterization of the  
 125 DRIR’s mixing time using a measure of the sound field’s diffuseness, before showing that  
 126 re-synthesizing the reverberation tail as a zero-mean Gaussian noise in the SHD preserves  
 127 the late field’s spatial properties.

### 128 1. EDR analysis

129 The EDR is a time-frequency extension of Schroeder’s reverse-integrated broadband en-  
 130 ergy decay curve (EDC), from which frequency-dependent decay envelope parameters can be  
 131 extracted by analyzing each frequency bin individually (Jot *et al.*, 1997). We begin our anal-  
 132 ysis by identifying the exponential decay section of the reverse-integrated curve presented

133 by the EDR at each frequency bin. In dB scale (such that exponential sections become  
 134 linear), this curve is first segmented using an adaptive Ramer-Douglas-Peucker (RDP) algo-  
 135 rithm (Prasad *et al.*, 2012) in order to help identify the different sections (early reflections,  
 136 exponential decay, and noise floor).

137 The noise floor limit point  $\{P_{\text{noise}}, t_{\text{lim}}\}$  can be found by fitting the theoretical dB-scale  
 138 profile of a reverse-integrated constant-power noise to the curve segments (see the shaded  
 139 area on Fig. 2). An additional headroom above this noise profile is then adaptively deter-  
 140 mined (see below) to ensure the limiting point  $\{P_{\text{noise}}, t_{\text{lim}}\}$  belongs to the exponential decay  
 141 section of the curve, thereby avoiding discontinuities when prolonging the reverberation en-  
 142 velope for tail re-synthesis. Finally, any non-exponentially decaying early reflection regimes  
 143 are discarded by selecting an appropriate starting segmentation point ( $t_{\text{start}}$ , see Fig. 2) using  
 144 a criterion on the local slopes of the curve segments up until  $t_{\text{lim}}$  (early segments to discard  
 145 are assumed to be shorter and have significantly different local slopes than those belonging  
 146 to exponential decays). The exponential decay section is thus delimited by  $t_{\text{start}}$  and  $t_{\text{lim}}$   
 147 and the reverberation time ( $T_{60}$ ) and initial power ( $P_0$ ) values can be determined by fitting  
 148 an ideal decay envelope model.

149 In the case of a single-slope decay, the envelope parameters can be found by performing  
 150 a linear regression on the identified decay section of the dB-scale curve. For multiple-slope  
 151 decays, such as those observed in certain configurations of coupled volumes (Cremer *et al.*,  
 152 1982), a parameter-space search can be performed in order to fit the model to the measured

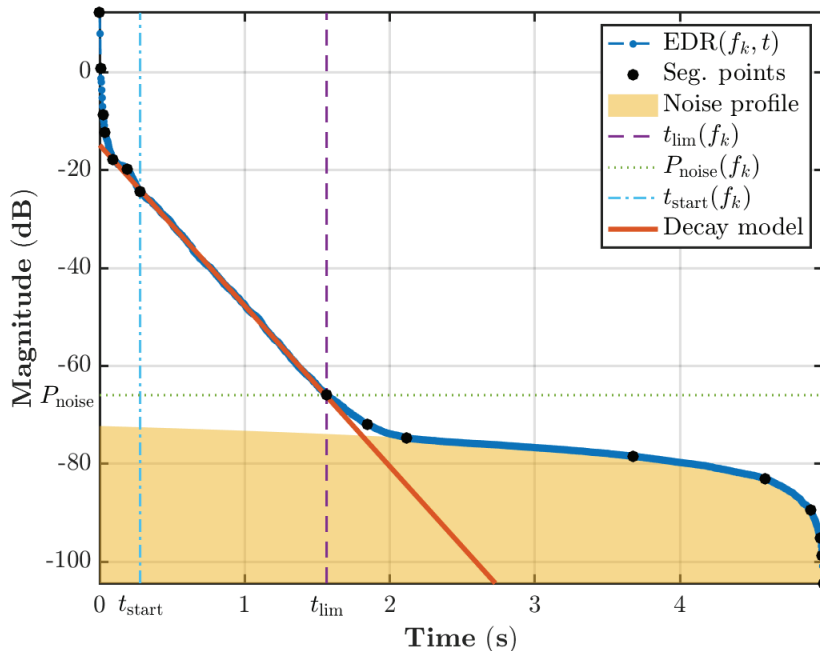


FIG. 2. (Color online.) EDR analysis schematic for a given frequency bin. The reverse-integrated decay curve is first segmented (black points). The noise floor (shaded area) is then identified, along with the noise floor limiting point  $\{P_{\text{noise}}, t_{\text{lim}}\}$  (dotted and dashed lines, respectively). Early decay sections are avoided by identifying  $t_{\text{start}}$  (dash-dot line), and the exponential decay model is fitted between  $t_{\text{start}}$  and  $t_{\text{lim}}$ .

153 decay (Xiang *et al.*, 2011). In general, if we consider the global energy envelope of a system  
 154 of  $C$  coupled volumes to be a sum of  $C$  exponential decays:

$$\text{ENV}(f, t, \mathbf{\Lambda}) = \sum_{i=1}^C P_{0,i}(f) e^{-2\delta_i(f)t}, \quad (3)$$

155 where  $\delta_i(f)$  are the frequency-dependent decay coefficients, related to the  $T_{60}(f)$  by  
 156  $T_{60}(f) = 3 \ln(10)/\delta(f)$ , and  $\mathbf{\Lambda}$  denotes the parameter vector containing the  $P_{0,i}$  and  $\delta_i$   
 157 values, then the ideal integrated decay curve is given by (see also [Jot et al., 1997](#))

$$\begin{aligned}
 \widehat{\text{EDR}}(f, t, \mathbf{\Lambda}) &= \int_t^\infty \text{ENV}(f, \tau, \mathbf{\Lambda}) d\tau \\
 &= \sum_{i=1}^C \frac{P_{0,i}(f)}{2\delta_i(f)} e^{-2\delta_i(f)t}.
 \end{aligned} \tag{4}$$

158 A model error  $\epsilon_{\text{mod}}$  can be defined as a simple mean-squared error (MSE):

$$\epsilon_{\text{mod}}(f) = \frac{1}{N_{\text{fit}}} \sqrt{\sum_{n=n_s}^{n_e} \left[ \text{EDR}_{\text{dB}}(f, t_n) - \widehat{\text{EDR}}_{\text{dB}}(f, t_n, \mathbf{\Lambda}) \right]^2}, \tag{5}$$

159 where  $N_{\text{fit}} = n_e - n_s + 1$ , with  $n_s$  the discrete time index such that  $t_{n_s} = t_{\text{start}}$  and  
 160 similarly  $n_e$  such that  $t_{n_e} = t_{\text{lim}}$ . This error can then be used as a loss function (or inversely  
 161 as a likelihood) in order to perform the parameter search using an expectation-maximisation  
 162 (EM) or maximum-likelihood (ML) algorithm. At each frequency bin, the parameter space is  
 163 of dimension  $2C$ , since for each exponential decay both  $P_{0,i}(f)$  and  $\delta_i(f)$  must be estimated.  
 164 To optimize the EM and avoid the detection of false local likelihood maxima, the algorithm  
 165 is initialized using linear regressions performed on EDC segments defined by re-applying the  
 166 adaptive RDP algorithm between  $t_{\text{start}}$  and  $t_{\text{lim}}$ .

167 The model error can additionally be used to adjust the headroom above the fitted ideal  
 168 noise profile mentioned above. The procedure described above (segmentation, noise fitting,  
 169 start point detection, and decay parameter search) is reiterated for several headroom values,  
 170 and the result with the highest overall likelihood (lowest error) is chosen. The likelihood  
 171 function used in this work is based on the Akaike Information Criterion (AIC) ([Akaike](#),

172 1974) and can be written  $\mathcal{L} = 2 \log(1/\epsilon_{\text{mod}}) - 2C + \log(N_{\text{fit}})$ , where again  $C$  is the number of  
 173 coupled decays, and  $\log(N_{\text{fit}})$  is a regularization term used to promote fits made over longer  
 174 decay sections (i.e. for two fits with equal likelihood, the one made over a longer section of  
 175 the EDR bin will be preferred).

## 176 2. *Diffuseness analysis and mixing time estimation*

177 As mentioned in section IB, replacing the non-decaying noise floor with a reverberation  
 178 tail synthesized as an exponentially-decaying zero-mean Gaussian noise assumes that the  
 179 late sound field described by the impulse response is fully diffuse. This leads to the classic  
 180 time-frequency limits for stochastic modelling of room reverberation, respectively the mixing  
 181 time and Schroeder frequency (Polack, 1988). The exploration of strategies for denoising in  
 182 the modal domain below the Schroeder frequency is left to future work; in this paper we will  
 183 apply the tail re-synthesis process across all frequencies, and note that for most reverberant  
 184 spaces the Schroeder frequency is low enough that the human auditory system is largely  
 185 insensitive to the modal reverberation below it. (This can be seen by comparing Schroeder’s  
 186 measure  $f_{\text{Sch}} \approx 2000\sqrt{\bar{T}_{60}/V}$ , where  $\bar{T}_{60}$  is a broadband measure of the reverberation time  
 187 and  $V$  is the volume of the space (Schroeder and Kuttruff, 1962), to equal-loudness contours  
 188 such as those given by the ISO226:2003 standard.)

189 Defining the mixing time, however, is crucial to the present work. Considering the afore-  
 190 mentioned requirement of a fully diffuse late sound field for synthesizing the prolongation  
 191 of the reverberation tail using a zero-mean Gaussian noise, we propose using a measure of  
 192 the sound field’s diffuseness in order to estimate the moment the DRIR becomes maximally

193 diffuse. Furthermore, in the following section [II B 3](#) we will show that re-synthesizing the  
194 late reverberation tail in the SHD guarantees that the resulting sound field will preserve  
195 these diffuseness properties.

196 Several measures of diffuseness have been proposed that directly exploit various charac-  
197 teristics of the SHD. The DirAc measure ([Ahonen and Pulkki, 2009](#)) uses the zeroth- and  
198 first-order components to define a sound intensity vector and analyze its temporal variation.  
199 [Jarrett et al. \(2012\)](#) use SHD inter-component coherence to define a “signal-to-diffuse ra-  
200 tio” (SDR) that is evaluated with respect to a directional signal with a given direction of  
201 arrival (DOA). Finally, the COMEDIE measure ([Epain and Jin, 2016](#)) exploits the eigen-  
202 decomposition of the SHD signal covariance matrix, which will approach the identity matrix  
203 in the case of a fully diffuse field. The COMEDIE measure was chosen for this work due  
204 to its increase in accuracy with SHD order (whereas the DirAc measure is limited to first-  
205 order signals), its relatively lightweight implementation, and its independence from external  
206 analyses (whereas the SDR measure relies on the DOA).

207 In typical “large” mixing spaces, diffuseness profiles tend to quickly reach a stable max-  
208 imum, as shown in [Fig. 3](#) for the Kraftzentrale event venue in Duisburg, Germany (an  
209 industrial-era factory hall approximately 84000 m<sup>3</sup> in volume). Estimating the mixing time  
210 then corresponds to identifying the moment the DRIR reaches its maximum diffuseness.  
211 The idea here is to first characterize the maximum diffuseness and then find when the DRIR  
212 reaches this maximum in a definitive manner after an initial period of instability due to

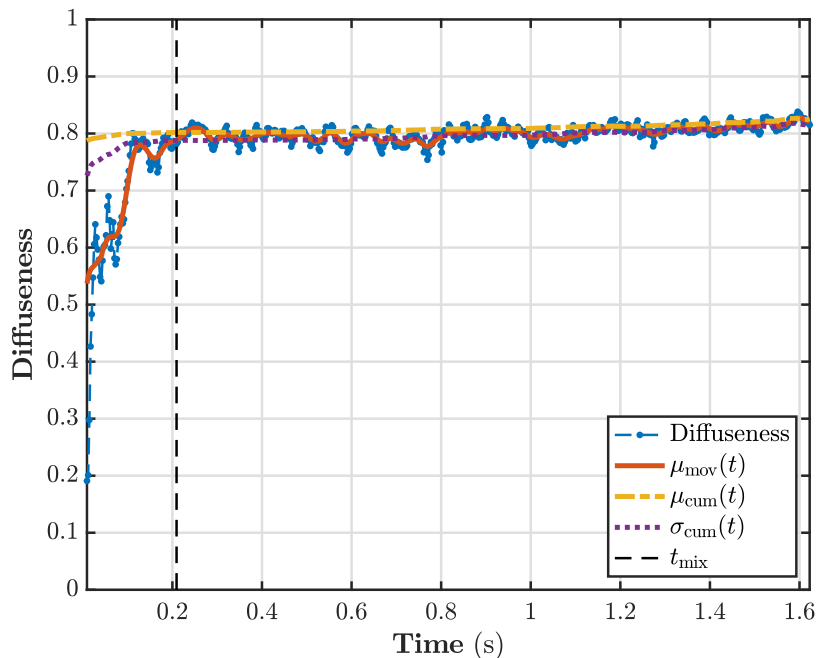


FIG. 3. (Color online.) COMEDIE diffuseness analysis (Epain and Jin, 2016) and mixing time estimation for a 4<sup>th</sup>-order SHD DRIR measured at the Kraftzentrale event venue in Duisburg, Germany, using an mh acoustics Eigenmike<sup>®</sup>. The calculated diffuseness curve is smoothed using a Gaussian kernel moving average ( $\mu_{\text{mov}}$ ). An inverse cumulative average ( $\mu_{\text{cum}}$ ) and standard deviation ( $\sigma_{\text{cum}}$ ) are further used to identify the onset of the maximum diffuseness and thereby estimate the mixing time ( $t_{\text{mix}}$ ).

213 coherent early reflections. This can be done by means of an appropriately-sized moving  
 214 average,

$$\mu_{\text{mov}}(t_i) = \sum_{n=i}^{i+N_w-1} w(t_n - t_i)d(t_n), \quad (6)$$

215 a reverse-cumulative average,

$$\mu_{\text{cum}}(t_i) = \frac{1}{N_d - i + 1} \sum_{n=i}^{N_d} d(t_n), \quad (7)$$

216 and a reverse-cumulative standard deviation,

$$\sigma_{\text{cum}}(t_i) = \sqrt{\frac{1}{N_d - i} \sum_{n=i}^{N_d} [d(t_n) - \mu_{\text{cum}}(t_n)]^2}, \quad (8)$$

217 where  $i = 1, 2, \dots, (N_d - N_w + 1)$  with  $N_d$  the length of the diffuseness data  $d(t_i)$  and  
 218  $w$  a chosen averaging kernel of length  $N_w$ . In this work (see Fig. 3), a 24-point Gaussian  
 219 kernel was used to calculate  $\mu_{\text{mov}}$  on diffuseness data obtained using a 1024-sample, 87.5%  
 220 overlapping short-term Fourier transform and mathematical expectations estimated by a  
 221 subsequent 8-frame average (at a 48 kHz sampling rate, this corresponds to a 40.0 ms  
 222 average for diffuseness points and a 101 ms total average for  $\mu_{\text{cov}}$ ). The mixing time is then  
 223 determined by

$$t_{\text{mix}} = \min(t_{\text{diff}}), \quad (9)$$

224 where the time values  $t_{\text{diff}}$  satisfy

$$\sqrt{[\mu_{\text{mov}}(t_{\text{diff}}) - \mu_{\text{cum}}(t_{\text{diff}})]^2} \leq \sigma_{\text{cum}}(t_{\text{diff}}). \quad (10)$$

225 Additional checks can subsequently be performed to ensure that no  $\mu_{\text{mov}}$  values are below  
 226 a certain threshold from  $\mu_{\text{cum}}$  after this time (e.g. corresponding to late-arriving discrete  
 227 echoes), adjusting  $t_{\text{mix}}$  to a satisfying  $t_{\text{diff}}$  value if necessary. Further validation tests on the



228 value of the maximum diffuseness may also be included (e.g. a diffuseness maximum below  
 229 0.5 may not be considered “maximally diffuse”).

230 We now need to define a condition for re-synthesizing the reverberation tail using a  
 231 zero-mean Gaussian noise: if the DRIR reaches its mixing time before decaying below the  
 232 noise floor, the stochastic model can be used as first proposed by [Jot \*et al.\* \(1997\)](#). However,  
 233 whereas the mixing time is a broadband property, the EDR analysis described above returns  
 234 a frequency-dependent noise floor limiting time  $t_{\text{lim}}(f)$ . To get a global value for the noise  
 235 floor limiting time, we use the  $t_{\text{lim}}(f)$  values determined for the SHD-encoded DRIR’s  $Y_{0,0}$   
 236 (omnidirectional) component and perform a perceptually-weighted average over the audible  
 237 frequency range. This average is weighted according to the ITU-R 468 standard noise filter  
 238 and then evaluated over Bark-scale frequency bands in order to avoid the over-weighting of  
 239 higher-frequency bins due to the linear frequency scale of the Fourier transform.

240 We denote the resulting value  $\bar{t}_{\text{lim}}$ , and the condition can then be written  $t_{\text{mix}} < \bar{t}_{\text{lim}}$ .  
 241 If it is verified, tail re-synthesis may be performed using a zero-mean Gaussian noise as  
 242 described below, with the perceptual considerations above ensuring that any  $t_{\text{lim}}(f)$  values  
 243 smaller than  $t_{\text{mix}}$  should have a limited perceptual impact (future work is planned to further  
 244 strengthen this aspect, e.g. by taking into account the corresponding  $P_{\text{noise}}(f)$  values). If the  
 245 condition is not verified, however, alternative methods of noise reduction must be considered  
 246 (see the conclusion in section [IV](#) below).

247 **3. Diffuse tail synthesis**

248 We now show that re-synthesizing the reverberation tail as a zero-mean Gaussian noise  
 249 in the SHD preserves the spatial incoherence properties of the late reverberation field. In  
 250 the SHD, the signal measured by a SMA in the presence of a perfectly diffuse field is of the  
 251 form

$$X_{l,m}^{\text{diff}}(f, t) = \sqrt{P_{\text{diff}}(f, t)} b_l(f) \int_{\Omega \in S^2} \Phi(f, \Omega, t) Y_{l,m}(\Omega) d\Omega, \quad (11)$$

252 where  $P_{\text{diff}}(f, t)$  is the diffuse field power envelope,  $\Phi(f, \Omega, t) = e^{i\varphi(f, \Omega, t)}$  with  $\varphi(f, \Omega, t)$   
 253 the independent and uncorrelated plane wave phase such that  $|\Phi(f, \Omega, t)| = 1 \forall f, \Omega, t$  and  
 254  $E\{\Phi(f, \Omega, t)\Phi^*(f, \Omega', t)\} = \delta_{\Omega, \Omega'}$  (with  $\delta$  representing the Kronecker delta and  $E\{\cdot\}$  math-  
 255 ematical expectation), and  $b_l(f)$  are the aforementioned array mode strengths (or holo-  
 256 graphic functions). It can be shown that this leads to a spatial coherence of  $\gamma_{l,m;l',m'}^{\text{diff}}(f, t) =$   
 257  $0 \forall (l, m) \neq (l', m')$  (Jarrett *et al.*, 2012) due to the orthogonality of the spherical harmonics  
 258 and the spatial independence of the plane wave phases.

259 On the other hand, synthesizing a zero-mean Gaussian noise of power  $P_{l,m}^{\text{diff}}(f, t)$  and  
 260 random phase  $\Phi_{l,m}(f, t) = e^{i\varphi_{l,m}(f, t)}$  per SHD component gives a cross-power spectral density  
 261 (PSD) of

$$\begin{aligned} \hat{\Psi}_{l,m;l',m'}^{\text{diff}}(f, t) &= E\left\{\hat{X}_{l,m}^{\text{diff}}(f, t)\hat{X}_{l',m'}^{\text{diff}*}(f, t)\right\} \\ &= P_{l,m}^{\text{diff}}(f, t)P_{l',m'}^{\text{diff}}(f, t)\delta_{l,m;l',m'}, \end{aligned} \quad (12)$$

262 and therefore the same diffuse field spatial coherence:

$$\begin{aligned}
 \hat{\gamma}_{l,m;l',m'}^{\text{diff}}(f, t) &= \frac{\hat{\Psi}_{l,m;l',m'}^{\text{diff}}(f, t)}{\sqrt{\hat{\Psi}_{l,m;l,m}^{\text{diff}}(f, t)}\sqrt{\hat{\Psi}_{l',m';l',m'}^{\text{diff}}(f, t)}} \\
 &= 0 \quad \forall (l, m) \neq (l', m').
 \end{aligned}
 \tag{13}$$

263 It can also be shown that synthesizing a zero-mean Gaussian noise per SHD component  
 264 leads to an SHD covariance matrix that approaches the identity matrix in the same way  
 265 as a diffuse field of  $N \gg (L + 1)^2$  independent and uncorrelated plane waves, as originally  
 266 demonstrated by [Epain and Jin \(2016\)](#) for the COMEDIE diffuseness measure, provided a  
 267 normalized covariance calculation is used (thereby imitating a coherence). Although the use  
 268 of individual power envelopes per SHD component does not guarantee an ideally diffuse field,  
 269 it does guarantee at least a fully incoherent field, and is furthermore necessary to account for  
 270 both the order-dependent frequency response of the SHD components ([Daniel and Moreau,](#)  
 271 [2004](#)) and any deviations from perfect isotropy which could introduce continuity artefacts  
 272 at  $\{P_{\text{noise}}, t_{\text{lim}}\}$  when prolonging the reverberation tail.

### 273 C. Summary of denoising process

274 The full denoising process (outlined in Fig. 1) can thus be summarized as follows:

- 275 1. *Measurement artefact reduction.* The procedure described in section [II A](#) is applied to  
 276 the raw ESM recording data of each SMA microphone channel.
- 277 2. *Inverse-sweep convolution and SHD transform.* The resulting “cleaned” ESM mea-  
 278 surement is convolved with a time-reversed and amplitude-corrected version of the

279 excitation sweep signal as per Farina (2000) to obtain an IR for each microphone  
 280 channel. This multi-channel IR is then transformed to the SHD according to the  
 281 theory outlined in section IA.

282 3. *Mixing time analysis.* Diffuseness analysis is performed in the SHD, leading to an  
 283 estimation of the mixing time as presented in section IIB2.

284 4. *EDR analysis and validation of diffuse field hypothesis.* EDR analysis is performed per  
 285 SHD component in order to extract the reverberation tail decay envelope parameters  
 286 ( $T_{60}(f)$  and  $P_0(f)$ ) and noise floor limit points  $\{P_{\text{noise}}, t_{\text{lim}}\}(f)$ . The  $t_{\text{lim}}(f)$  values ob-  
 287 tained for the omnidirectional  $Y_{0,0}$  component are averaged over the audible frequency  
 288 range in order to estimate the broadband noise floor limiting time and confirm (or  
 289 invalidate) the diffuse field hypothesis required for tail re-synthesis using a zero-mean  
 290 Gaussian noise.

291 5. *Tail re-synthesis.* The late reverberation tail is re-synthesized using a zero-mean Gaus-  
 292 sian noise per SHD component, which preserves spatial incoherence as shown above.  
 293 For every SHD component channel, each frequency bin of the re-synthesized tail is  
 294 made to decay according to the corresponding parameters extracted from the DRIR,  
 295 and is then used to replace the corresponding DRIR frequency bin starting at  $t_{\text{lim}}(f)$ .

### 296 III. APPLICATION TO MEASURED DRIR

297 In this section we show the effects of applying the denoising process described above to  
 298 DRIRs measured in various locations and conditions. A qualitative overview of the results

299 is first presented, followed by a brief discussion of methods leading to a more quantitative  
 300 assessment of the procedure’s performance.

### 301 **A. Measurement artefact reduction**

302 Figs. 4 and 5 illustrate the application of the artefact reduction method described in  
 303 section II to a single microphone channel of an ESM measurement performed at the Chris-  
 304 tuskirche in Karlsruhe, Germany (a late 19th-century church with a large open dome-like  
 305 nave). Fig. 4 (a) shows several impulsive artefacts occurring over the course of the ESM  
 306 measurement signal (averaged over four repetitions), while Fig. 5 (a) illustrates how these  
 307 turn into repeated inverse-sweep artefacts when the ESM measurement signal is convolved  
 308 with the time-reversed and amplitude-corrected excitation signal as per Farina (2000).  
 309 Figs. 4 (b) and 5 (b) show the effect of the artefact reduction procedure on the ESM mea-  
 310 surement signal and resulting IR, respectively. Finally, Figs. 4 (c) and 5 (c) highlight the  
 311 time-frequency points identified as artefacts as well as their magnitude differences before  
 312 and after reduction.

313 The spectrograms shown in Figs. 4 and 5 are obtained by performing a moving time  
 314 average over 8 frames of short-term Fourier transform magnitudes (with 87.5% overlapping  
 315 frames of 1024 samples at a 48 kHz sampling rate, this corresponds to a total averaging  
 316 length of 40 ms).

317 The removal of the inverse-sweep-type artefacts revealed in Fig. 5 is crucial in ensuring  
 318 that the reverse-integration of the IR’s noise floor approaches the theoretical profile fitted  
 319 to identify the noise floor limit point  $\{P_{\text{noise}}, t_{\text{lim}}\}(f)$ , as in Fig. 2 (see section II B 1). To

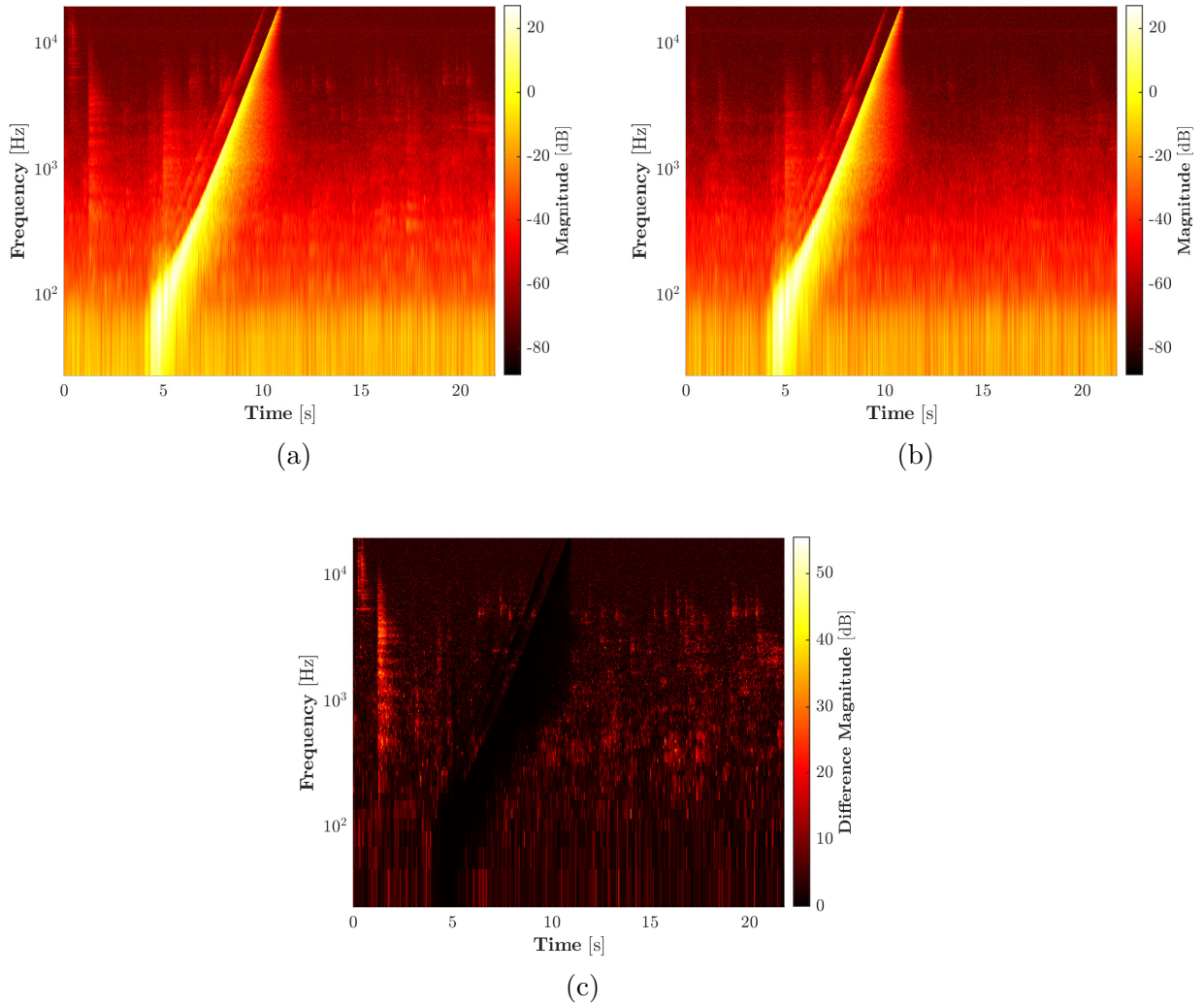


FIG. 4. (Color online.) Artefact reduction applied to a single microphone channel of an ESM measurement performed at the Christuskirche in Karlsruhe, Germany, using an mh acoustics Eigenmike<sup>®</sup>. (a) Spectrogram of the raw ESM measurement signal (averaged over four repetitions), with several impulsive sounds present. (b) Spectrogram of the ESM measurement signal after artefact reduction. (c) Spectrogram difference between (a) and (b).

320 further illustrate this, Fig. 6 compares the EDR profile from the Christuskirche DRIR's  $Y_{0,0}$   
 321 component for one frequency (2461 Hz) before and after application of the artefact reduction  
 322 process.

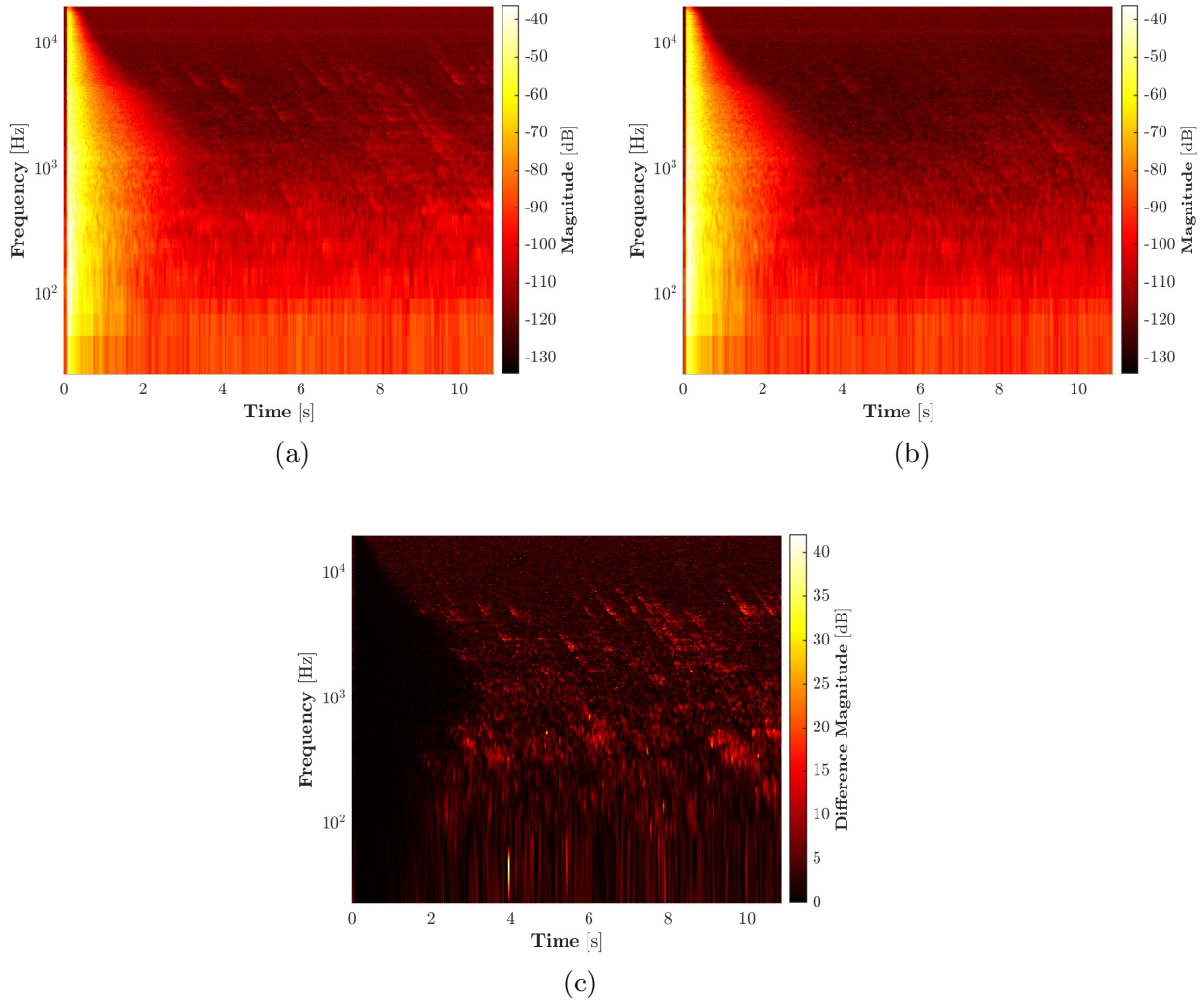


FIG. 5. (Color online.) Artefact reduction applied to a single microphone channel of an ESM measurement performed at the Christuskirche in Karlsruhe, Germany, using an mh acoustics Eigenmike<sup>®</sup>. (a) Spectrogram of the original IR, after inverse-sweep convolution with the raw ESM measurement signal (without artefact reduction). (b) Spectrogram of the IR obtained by inverse-sweep convolution with the artefact-reduced ESM measurement signal. (c) Spectrogram difference between (a) and (b).

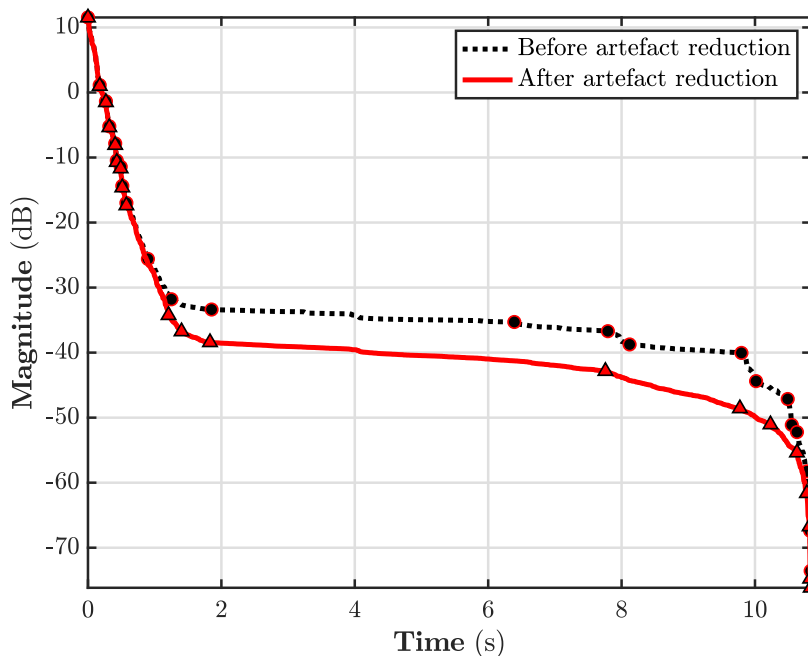


FIG. 6. (Color online.) EDR profile of the Christuskirche DRIR's omnidirectional  $Y_{0,0}$  component for one frequency (2461 Hz); before (black dashed line) and after (red solid line) artefact reduction. Circle and triangle markers represent adaptive Ramer-Douglas-Peucker segmentation points (see section II B 1).

323 Finally, in an attempt to quantify the amount of artefact reduction, we define an *artefact-*  
 324 *to-total-energy ratio* as the total artefact energy (i.e. the energy of spectral outliers according  
 325 to the definition given in section II A) versus the total signal energy in a given frame:

$$\eta(t) = \frac{\sum_{k=0}^K |\tilde{X}(f_k, t)|^2}{\sum_{k=0}^K |X(f_k, t)|^2},$$

$$\tilde{X}(f_k, t) = \begin{cases} X(f_k, t), & |X(f_k, t)| > \xi(f, t) \\ 0 & \end{cases} . \quad (14)$$



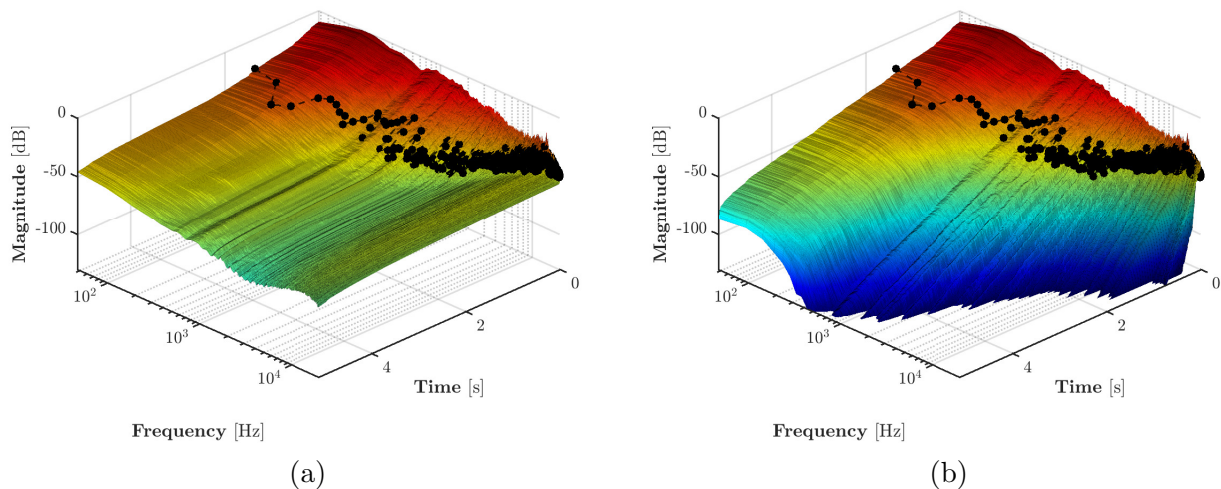


FIG. 7. (Color online.) EDRs of the Kraftzentrale DRIR’s omnidirectional  $Y_{0,0}$  component, (a) before and (b) after reverberation tail re-synthesis. The black dotted line shows the  $t_{\text{lim}}$  value for each frequency bin.

326 Thus  $\eta(t) = 0$  means that no outliers were found in the time frame, whereas  $\eta(t) = 1$   
 327 corresponds to an entirely outlying time frame. In the current example (the Christuskirche  
 328 DRIR), this measure averaged to  $\bar{\eta} = 0.274$  over the four sweep repetitions.

### 329 B. Reverberation tail re-synthesis

330 Figure 7 illustrates the effect of the tail re-synthesis procedure on the EDR of the om-  
 331 nidirectional  $Y_{0,0}$  component of the Kraftzentrale DRIR. The arbitrary dynamic range for  
 332 synthesis is chosen to match that of the signal bit depth (193 dB at 32 bits) at the most  
 333 perceptually important frequencies (again using the ITU-R 468 standard), although Fig. 7  
 334 is shown over 130 dB to match the depth of human hearing.

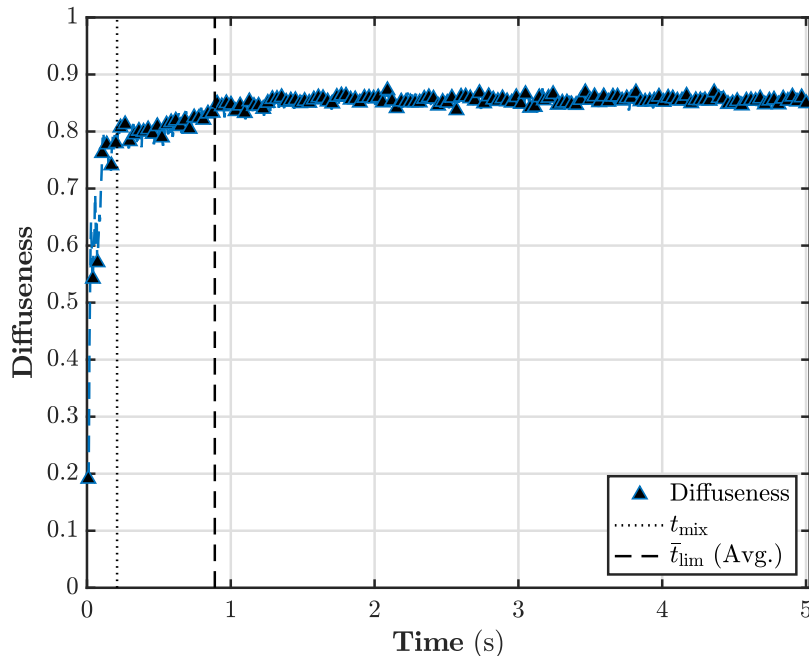


FIG. 8. (Color online.) COMEDIE diffuseness for the Kraftzentrale DRIR, after tail re-synthesis.

The  $t_{\text{mix}}$  and average  $\bar{t}_{\text{lim}}$  values (dotted and dashed lines, respectively) are shown as temporal references.

335 As mentioned throughout this paper, the crucial condition for successfully denoising  
 336 DRIRs by reverberation tail re-synthesis is that the late field’s diffuseness properties must  
 337 be preserved. To confirm that the proposed denoising procedure achieves this, Fig. 8 shows  
 338 the COMEDIE diffuseness profile for the Kraftzentrale DRIR: the diffuseness maximum  
 339 reached at  $t_{\text{lim}}$  (dotted line) is successfully extended and maintained beyond the average  
 340  $\bar{t}_{\text{lim}}$  (dashed line). Note that the COMEDIE diffuseness increases slightly from  $t_{\text{mix}}$  to  $\bar{t}_{\text{lim}}$ ,  
 341 which may be due to the method’s additional sensitivity to ideally diffuse signals versus  
 342 large numbers of plane waves, as initially noted by [Epain and Jin \(2016\)](#).

343 **IV. CONCLUSION**

344 This paper has addressed the problem of removing the non-decaying noise floor inevitably  
345 present in DRIRs measured with SMAs and replacing it with a valid extension of the  
346 exponentially-decaying late reverberation tail. Building on previous research showing that  
347 this is possible for so-called “mixing” spaces by synthesizing the late reverberation as a  
348 zero-mean Gaussian noise and parameterizing its decay envelope by analyzing the EDR, we  
349 have demonstrated that performing this synthesis in the SHD guarantees preservation of  
350 the late field’s spatial incoherence. Additionally, we have shown that including an artefact  
351 reduction step before inverse-sweep convolution of the ESM measurement signal improves  
352 identification of the noise floor during EDR analysis. As a collateral development, we have  
353 also proposed an estimate of the mixing time using measures of DRIR diffuseness in the  
354 SHD.

355 Further work on this topic can be organized around three main themes. First, the question  
356 of appropriately determining the number of coupled decays to consider in multi-slope cases  
357 must be addressed to avoid over-fitting and ensure that the detected model satisfies coupled-  
358 volume theory. Second, cases where the late reverberation field presents highly anisotropic  
359 energy distributions must be further investigated, as the spatial symmetry of the SHD will  
360 not enable proper re-synthesis using the method presented in this paper. Finally, techniques  
361 must be developed for cases where the reverberation tail cannot be considered diffuse before  
362 reaching the noise floor, i.e. spaces that cannot be considered traditionally “mixing” and  
363 whose late reverberation cannot be modeled as a stochastic process.

364 **ACKNOWLEDGMENTS**

365 This work was funded in part by a doctoral research grant from the École doctorale  
366 Informatique, Télécommunications, et Électronique (Paris) at Sorbonne Université. The  
367 authors would additionally like to thank Augustin Muller (IRCAM) and Pedro Garcia-  
368 Velazquez (Le Balcon) for their extensive DRIR measurements, as well as Franck Zagala  
369 (PhD candidate, Sorbonne Université/IRCAM) for having provided the foundations of the  
370 EDR analysis and diffuse reverberation tail re-synthesis algorithms.

371

372 Ahonen, J., and Pulkki, V. (2009). “Diffuseness Estimation Using Temporal Variation of  
373 Intensity Vectors,” in *Proceedings of the 2009 IEEE Workshop on Applications of Signal  
374 Processing to Audio and Acoustics*, New Paltz, U.S.A., pp. 285–288, doi: [10.1109/ASPAA.  
375 2009.5346496](https://doi.org/10.1109/ASPAA.2009.5346496).

376 Akaike, H. (1974). “A New Look at the Statistical Model Identification,” *IEEE Transactions  
377 on Automatic Control* **AC-19**(6), 716–723, doi: [10.1109/TAC.1974.1100705](https://doi.org/10.1109/TAC.1974.1100705).

378 Cabrera, D., Lee, D., Yadav, M., and Martens, W. L. (2011). “Decay Envelope Manipulation  
379 of Room Impulse Responses: Techniques for Auralization and Sonification,” in *Proceedings  
380 of Acoustics '11*, Gold Coast, Australia, pp. 52–56.

381 Carpentier, T., Szpruch, T., Noisternig, M., and Warusfel, O. (2013). “Parametric Control  
382 of Convolution-Based Room Simulators,” in *Proceedings of the 2013 International Sympo-  
383 sium on Room Acoustics*, Toronto, Canada.

- 384 Cremer, L., Müller, H. A., and Schultz, T. J. (1982). *Principles and Applications of Room*  
385 *Acoustics, vol. 1* (Applied Science Publishers, Barking, England).
- 386 Daniel, J., and Moreau, S. (2004). “Further Study of Sound Field Coding with Higher Order  
387 Ambisonics,” in *Proceedings of the 116<sup>th</sup> Audio Engineering Society Convention*, Berlin,  
388 Germany, <http://www.aes.org/e-lib/browse.cfm?elib=12789>.
- 389 Driscoll, J. R., and Healy, D. M. J. (1994). “Computing Fourier Transforms and Convo-  
390 lutions on the 2-Sphere,” *Advances in Applied Mathematics* **15**, 202–250, doi: [10.1006/](https://doi.org/10.1006/aama.1994.1008)  
391 [aama.1994.1008](https://doi.org/10.1006/aama.1994.1008).
- 392 Epain, N., and Jin, C. T. (2016). “Spherical Harmonic Signal Covariance and Sound  
393 Field Diffuseness,” *IEEE/ACM Transactions on Audio, Speech, and Language Process-*  
394 *ing* **24**(10), 1796–1807, doi: [10.1109/TASLP.2016.2585862](https://doi.org/10.1109/TASLP.2016.2585862).
- 395 Farina, A. (2000). “Simultaneous Measurement of Impulse Response and Distortion with a  
396 Swept-Sine Technique,” in *Proceedings of the 108<sup>th</sup> Audio Engineering Society Convention*,  
397 Paris, France, <http://www.aes.org/e-lib/browse.cfm?elib=10211>.
- 398 Guski, M., and Vorländer, M. (2014). “Comparison of Noise Compensation Methods for  
399 Room Acoustic Impulse Response Evaluations,” *Acta Acustica United with Acustica*  
400 **100**(2), 320–327, doi: [10.3813/AAA.918711](https://doi.org/10.3813/AAA.918711).
- 401 ISO226:2003 (2003). “Acoustics – Normal Equal-Loudness-Level Contours” (Interna-  
402 tional Organization for Standardization, Geneva, Switzerland), [https://www.iso.org/](https://www.iso.org/standard/34222.html)  
403 [standard/34222.html](https://www.iso.org/standard/34222.html).
- 404 ISO8000-2:2009(E) (2009). “Quantities and Units – Part 2: Mathematical Signs and Sym-  
405 bols to be Used in the Natural Sciences and Technology” (International Organization for

- 406 Standardization, Geneva, Switzerland), <https://www.iso.org/standard/64973.html>.
- 407 Jarrett, D. P., Thiergart, O., Habets, E. A. P., and Naylor, P. A. (2012). “Coherence-  
408 Based Diffuseness Estimation in the Spherical Harmonic Domain,” in *Proceedings of the*  
409 *27<sup>th</sup> IEEE Convention of Electrical and Electronics Engineers in Israel*, Eilat, Israel, doi:  
410 [10.1109/EEEI.2012.6377148](https://doi.org/10.1109/EEEI.2012.6377148).
- 411 Jot, J.-M., Cerveau, L., and Warusfel, O. (1997). “Analysis and Synthesis of Room Re-  
412 verberation Based on a Statistical Time-Frequency Model,” in *Proceedings of the 103<sup>rd</sup>*  
413 *Audio Engineering Society Convention*, New York, U.S.A., [http://www.aes.org/e-lib/](http://www.aes.org/e-lib/browse.cfm?elib=7150)  
414 [browse.cfm?elib=7150](http://www.aes.org/e-lib/browse.cfm?elib=7150).
- 415 Noisternig, M., Carpentier, T., Szpruch, T., and Warusfel, O. (2014). “Denoising of Direc-  
416 tional Room Impulse Responses Measured with Spherical Microphone Arrays,” in *Proceed-*  
417 *ings of the 40<sup>th</sup> Annual German Congress on Acoustics (DAGA)*, Oldenburg, Germany, pp.  
418 600–601, [http://pub.dega-akustik.de/DAGA\\_2014/data/articles/000292.pdf](http://pub.dega-akustik.de/DAGA_2014/data/articles/000292.pdf).
- 419 Noisternig, M., Zotter, F., and Katz, B. F. G. (2011). “Reconstructing Sound Source Di-  
420 rectivity in Virtual Acoustic Environments,” in *Principles and Applications of Spatial*  
421 *Hearing*, edited by Y. Suzuki, D. Brungart, Y. Iwaya, K. Iida, D. Cabrera, and H. Kato  
422 (World Scientific Publishing Co. Pte. Ltd.), pp. 357–373.
- 423 Polack, J.-D. (1988). “La transmission de l’énergie sonore dans les salles,” Ph.D. thesis,  
424 Université du Maine.
- 425 Prasad, D. K., Leung, M. K., Quek, C., and Cho, S. Y. (2012). “A Novel Framework for  
426 Making Dominant Point Detection Methods Non-Parametric,” *Image and Vision Comput-*  
427 *ing* **30**(12), 843–859, doi: [10.1016/j.imavis.2012.06.010](https://doi.org/10.1016/j.imavis.2012.06.010).

- 428 Rafaely, B. (2005). “Analysis and Design of Spherical Microphone Arrays,” IEEE Transac-  
429 tions on Speech and Audio Processing **13**(1), 135–143, doi: [10.1109/TSA.2004.839244](https://doi.org/10.1109/TSA.2004.839244).
- 430 Schroeder, M. R. (1962). “Natural-Sounding Artificial Reverberation,” Journal of the Audio  
431 Engineering Society **10**(3), 219–223.
- 432 Schroeder, M. R., and Kuttruff, H. (1962). “On Frequency Response Curves in Rooms:  
433 Comparison of Experimental, Theoretical, and Monte Carlo Results for the Average Fre-  
434 quency Spacing between Maxima,” The Journal of the Acoustical Society of America **34**(1),  
435 76, doi: [10.1121/1.1909022](https://doi.org/10.1121/1.1909022).
- 436 Xiang, N., Goggans, P., Jasa, T., and Robinson, P. (2011). “Bayesian Characterization of  
437 Multiple-Slope Sound Energy Decays in Coupled-Volume Systems,” The Journal of the  
438 Acoustical Society of America **129**(2), 741–752, doi: [10.1121/1.3518773](https://doi.org/10.1121/1.3518773).