

# Blind identification of Ambisonic reduced room impulse response

Srđan Kitić and Jérôme Daniel

**Abstract**—Recently proposed *Generalized Time-domain Velocity Vector (GTVV)* is a generalization of relative room impulse response in spherical harmonic (*aka* Ambisonic) domain, that allows for blind estimation of early-echo parameters: the directions and relative delays of individual reflections. However, the derived closed-form expression of GTVV mandates few assumptions to hold, most important being that the impulse response of the reference signal needs to be a minimum-phase filter. In practice, the reference is obtained by spatial filtering towards the Direction-of-Arrival of the source, and the aforementioned condition is bounded by the performance of the applied beamformer (and thus, by the Ambisonic array order). In the present work, we circumvent this problem by directly modeling the impulse responses constituting the GTVV time series, which permits not only to relax the initial assumptions, but also to extract the information therein in a more consistent and efficient manner, entering the realm of blind system identification. Experiments using simulated and recorded room impulse responses confirm the effectiveness of the proposed approach.

**Index Terms**—blind identification, Ambisonic, microphone array, RTF, Prony

## I. INTRODUCTION

Room Impulse Response (RIR) can be thought of as an “acoustic fingerprint” of the surrounding environment [1], and its importance in spatial audio processing cannot be overstated. It encodes the information about acoustic multipath - reverberation, which inevitably affects all indoor audio recordings. Traditionally often seen as a nuisance, reverberation is known to degrade the results of localization algorithms [2], worsen automatic speech recognition (ASR) and intelligibility [3], and negatively impact sound source separation [4]. Nevertheless, a number of recent works has demonstrated that *early* reverberation actually have a potential to improve performances in various tasks. These “echo-enabled” methods exploit early reflections to boost performance of acoustic localization [5], [6], [7], source separation [8], [9], [10], speech and sound event recognition [11], [12], [13], but also to address some unconventional problems such as localization behind soundproof obstacles [14], [15], [16], inference of room geometry [17], [18], [19], distance estimation [20], [21], [22], identification of room acoustic parameters [23], [24], [25] and acoustic matching [26], [27].

While the availability of pre-recorded RIRs would be, therefore, very beneficial for many echo-enabled applications, such

procedure demands specific equipment and skills. Moreover, a particular RIR is dependent on the given acoustic conditions, hence in dynamic scenes (*e.g.*, when microphone or source are mobile) one would require repeated RIR measurements, which is clearly impractical. Thus, there is a growing interest for adaptive blind system identification (BSI) methods - subject to certain assumptions, these are capable of inferring RIRs (up to a common delay and scale [28]) using only recorded audio signals. However, classical adaptive BSI methods (*e.g.*, multichannel least mean squares (MCLMS) algorithm [29]), have notable drawbacks, such as their sensitivity to noise and incorrect model order [30].

Recently, we have investigated properties of the so-called Generalized Time-domain Velocity Vector (GTVV) [31], a generalization of the well-known relative room impulse response in spherical harmonic (SH) domain [32]. The main advantage of GTVV over classical relative RIR is due to its reference signal, which is obtained by beamforming towards the (approximate) Direction-of-Arrival (DoA) of a far field source. Assuming that the reference signal is dominated by direct propagation, the GTVV representation admits a closed-form expression that can be used to directly infer the DoAs and relative delays of individual wavefronts (including the direct component). To satisfy this assumption, beamformer needs to be sufficiently selective, which is acceptable for Higher-Order Ambisonic (HOA) [33] arrays, but becomes prohibitive when prevalent [34] First Order Ambisonic (FOA) arrays are used. Indeed, the beam width of FOA beamformers is too permissive [32], thus a number of non-attenuated reflections invalidates the former requirement.

The main contribution of this work is a method for the identification of Ambisonic RIRs with scale and delay shifting according to the principal wave front, directly from the GTVV imprint (thus, blindly) and regardless of the array order. Hence, we term this representation Reduced Room Impulse Response (RdRIR), all the more that we are primarily interested in extracting the early part of Ambisonic RIRs containing directional information, *i.e.* the “early echoes”. Our goal is also to retain low computational complexity, as well as to facilitate implementation. Therefore, we compare several algorithmic variants having different levels of complexity, and evaluate their estimation performance. We show through experiments on simulated and real impulse response data that the proposed methods are effective in extracting parameters of multiple wavefronts, under various acoustic conditions.

This work unifies and complements our previous contributions published as conference papers [7], [31]. The article is organized as follows: after a review of prior art given in Section II, we proceed to the signal model behind GTVV in

S. Kitić was with Orange Labs, France, at the time of writing this article. J. Daniel is with Orange Labs, France. The two authors have equally contributed to the present article.

This paper has supplementary downloadable material available at <http://ieeexplore.ieee.org>, provided by the authors. The material includes experimental results complementing those presented in the article. Contact [jerome.daniel@orange.com](mailto:jerome.daniel@orange.com) for further questions about this work.

Section III, and discuss its estimation and limitations. This is succeeded by presenting the RdRIR estimation methods in Section IV. The results of computer experiments are given in Section V. The article is concluded in Section VI.

**Notations:** Real- or complex-valued scalar variables are written in lowercase italic or greek alphabet, while we use boldface font for vectors (lowercase) and matrices (uppercase). Serif font is reserved for integers, with the uppercase serif denoting constant integer values. For a matrix  $M$ , we use  $m_{i,:}$  and  $m_{:,j}$  to specify its  $i^{\text{th}}$  row and  $j^{\text{th}}$  column, respectively, and  $m_{i,j}$  to denote the matrix entry at their intersection. The sets of natural, real and complex numbers are denoted by  $\mathbb{N}$ ,  $\mathbb{R}$  and  $\mathbb{C}$ , respectively. The Fourier transform (or the Discrete Fourier Transform, where appropriate) and its inverse are given by  $\mathcal{F}$  and  $\mathcal{F}^{-1}$ , while a variable in frequency domain is marked by the circumflex accent. The transpose and conjugate transpose operations are denoted by  $(\cdot)^{\text{T}}$  and  $(\cdot)^{\text{H}}$ , respectively. The notation  $x^{(i)}$  refers to the  $i^{\text{th}}$  iterate of some algorithmic variable  $x$ . For a vector-valued function  $\mathbf{x}(t)$ , its corresponding *delay-magnitude representation* is defined as  $\zeta_{\mathbf{x}}(t) = \|\mathbf{x}(t)\|_2$ .

## II. PRIOR ART

Research in BSI has flourished since the seminal work of Sato [35] on self-recovering equalization for digital communications. In the context of Single-Input-Multiple-Output (SIMO) systems, where the same source excites multiple channels (as in our problem setting), the concept of *cyclostationarity* in second order statistics (SOS) [36] was widely adopted. Different SOS variants have been proposed, based on channel cross-relation (CR) [29], subspace decomposition [37] and maximum likelihood estimation [38]. The CR technique has been particularly popular, and was later extended to adaptive BSI for acoustic channel identification, either in time [39], [40], or frequency domain [41], [42]. Nevertheless, while such methods have evolved in order to improve their robustness to noise, in general they are known to perform well only under sufficiently high Signal-to-Noise Ratio (SNR) [30]. Some CR variants have been tailored to estimation of early RIRs, the task referred to as “under-modeled BSI” in the literature [43], [44]. Contemporary approaches based on deep learning have only recently been employed for blind RIR estimation [45], [46], [13]. Some of these models achieve impressive performance on different metrics and datasets, but are currently limited to predicting only single channel RIRs. The interpolation of missing RIR channels of a circular microphone array, in the SIMO context, has been formulated and solved as an inverse problem regularized by deep prior in [47].

Relative Transfer Function (RTF) is a well-known concept in microphone array signal processing, that has been in widespread use for decades [48]. There are two distinct conveniences of RTF: i) it is obtained using only received multichannel signals as their ratio in frequency domain, and ii) it is theoretically a source signal-agnostic representation (thus, encoding only the propagating characteristics of the environment). Relative transfer functions have also been adopted in Ambisonic domain [32], [49], [50], yet sometimes under different names (*e.g.* relative harmonic coefficients (RHC) [51]

or frequency-domain velocity vector (FDVV) [7]). For FOA signals, its real part is aligned with *pseudointensity* vector [52], [53], a widely used alternative to steered beamforming for low-cost DoA estimation. The temporal representation of RTF is relative impulse response [54], [55] (again renamed to time-domain velocity vector (TDVV) in our earlier work [7]). An idea related to RTF and relative impulse response was discussed by Gölles and Zotter in [56], where they calculate the ratio of received Ambisonic signals directly in time domain.

Classical RTF and relative impulse response have been extended to generalized (frequency and time domain) velocity vectors [31], mentioned before and explained in detail in the next section<sup>1</sup>. The value of the beamformed reference has also been recognized in [58], [59], where the authors take it to be the proxy for the source signal, hence the obtained ratio is considered an approximation of the acoustic transfer function (ATF). However, this hypothesis can only be valid if the applied beamformer filters out all reflections, which is rarely the case. To alleviate this issue, in [59] the authors propose to use a time-frequency mask obtained by the improved direct-path-dominance test [60], which is nonetheless a costly estimator in terms of computational complexity. In [61], this representation has been used for denoising, under the assumption that the directions of acoustic reflections are known a priori.

As mentioned before, the central motivation of our work is extracting spatio-temporal information about the early echoes, and not identifying the complete propagation channels. A related work has been recently published by Shlomo and Rafaely [62], where they propose the phase aligned spatial correlation (PHALCOR) algorithm, for the same purpose. They obtain convincing results on simulated data, at the expense of high computational cost and somewhat intricate implementation, involving singular vector decomposition, sparse analysis and clustering.

## III. GENERALIZED VELOCITY VECTOR

In this section we recall and extend the concept of generalized velocity vector, introduced in [31]. In particular, generalized velocity vector definition is broadened to include frequency-dependent beamforming and attenuation, and its relation with RTF and pseudointensity is made more explicit. Moreover, we provide a closed-form expression for GTVV for the case where the dominant wavefront in the reference signal is an acoustic reflection (instead of direct sound).

### A. Signal model

Let  $\hat{\mathbf{b}}(f) \in \mathbb{C}^{(L+1)^2}$  denote the vector of concatenated spherical harmonic expansion coefficients (the “HOA channels”) up to order  $L$ , at frequency  $f$ . We assume that mode strength compensation [32] has been applied, and that the recorded signals are due to a far field point source at azimuth  $\theta_0$ ,

<sup>1</sup>Note that, in [57], Herzog and Habets have proposed another acoustic quantity termed *generalized intensity vector*, which, despite the naming similarity, is different from generalized velocity vectors discussed here.

elevation  $\phi_0$  and range  $d_0$  from the microphone array, in an indoor environment. We approximate  $\hat{\mathbf{b}}(f)$  as follows:

$$\hat{\mathbf{b}}(f) \approx \hat{s}(f)\hat{\chi}(f) \sum_{n=0}^{N-1} \hat{\nu}_n(f)\mathbf{y}_n e^{-j2\pi f\bar{\tau}_n} + \hat{\mathbf{e}}(f) \quad (1)$$

$$= \hat{\mathbf{x}}(f) + \hat{\mathbf{e}}(f). \quad (2)$$

In the expression above,  $\hat{s}(f)$  represents the source (excitation) wideband signal, such as speech, while  $\hat{\chi}(f)$  is an anti-aliasing filter applied before the analog-to-digital converter. The terms  $\hat{\nu}_n(f) \in (0, 1)$ ,  $\bar{\tau}_n \in \mathbb{R}^+$  and  $\mathbf{y}_n \in \mathbb{R}^{(L+1)^2}$  are the attenuation factor, Time-of-Arrival (ToA), and the real-valued SH encoding vector of the  $n^{\text{th}}$  acoustic wavefront, respectively. The plane wave expansion has been truncated to  $N$  wavefronts, aggregating direct propagation and early echoes in  $\hat{\mathbf{x}}(f)$ , while the late reverberation and diffuse noise are represented by an additive term  $\hat{\mathbf{e}}(f)$ . The ToA of the direct signal is  $\bar{\tau}_0 \approx d_0/c$  (where  $c$  is the speed of sound), while  $\bar{\tau}_{\max} = \max_n \bar{\tau}_n$  (roughly) corresponds to the mixing time [63] of the room. Note that there is an implicit dependency between  $N$ ,  $\bar{\tau}_{\max}$  and the geometric and acoustic properties of the environment.

The diffuse component  $\hat{\mathbf{e}}(f)$  is considered uncorrelated with the directional term  $\hat{\mathbf{x}}(f)$  [32]. The latter is modeled by the image source model (ISM) [64], which considers all reflections to be specular, and approximates the frequency-dependent factor  $\hat{\nu}(f)$  by absorption coefficient - an attenuation of the signal magnitude by a positive factor lower than 1. In the ISM model, hence, the phase shifts of individual wavefronts are only due to differences in lengths of their acoustic paths. This is a limitation of the model - in general,  $\hat{\nu}(f)$  is a complex variable that encodes the material absorption and air attenuation, and depends on the angle of incidence [1].

Given a beamformer  $\hat{\mathbf{w}}(f) \in \mathbb{C}^{(L+1)^2}$ , constrained by  $\hat{\mathbf{w}}(f)^H \mathbf{y}_0 = \beta_0 \in \mathbb{R}^+ = \text{const}$  (without loss of generality, we consider  $\beta_0 = 1$ ), *Generalized Frequency-domain Velocity Vector* (GFVV) has been defined [31] as

$$\hat{\mathbf{v}}(f) = \frac{\hat{\mathbf{x}}(f)}{\hat{\mathbf{w}}(f)^H \hat{\mathbf{x}}(f)} = \frac{\mathbf{y}_0 + \sum_{n=1}^{N-1} \hat{g}_n(f) e^{-j2\pi f\bar{\tau}_n} \mathbf{y}_n}{1 + \sum_{n=1}^{N-1} \hat{g}_n(f) \hat{\beta}_n(f) e^{-j2\pi f\bar{\tau}_n}}. \quad (3)$$

Here,  $\hat{g}_n(f) = \hat{\nu}_n(f)/\hat{\nu}_0(f) < 1$ ,  $\tau_n = \bar{\tau}_n - \bar{\tau}_0 > 0$  and  $\hat{\beta}_n(f) = \hat{\mathbf{w}}(f)^H \mathbf{y}_n / \beta_0$  denote the attenuation, delay and spatial response of the  $n^{\text{th}}$  reflected plane wave *relative* to the direct propagation component, respectively.

We further assume that  $\hat{\kappa}_n(f) = \hat{g}_n(f)\hat{\beta}_n(f)$  is sufficiently smooth, such that its time domain counterpart  $\kappa_n(t)$  is compact<sup>2</sup>. If  $\hat{\mathbf{w}}(f) := \mathbf{w}$  is a wideband beamformer, this condition is usually satisfied, since  $\hat{g}_n(f)$  (which can be thought of as a scaled attenuation factor), is often a slowly-varying function of frequency in standard rooms [1]. However, care should be taken with some data-dependent beamformers, such as Minimum Power Distortionless Response (MPDR), which may exhibit abrupt changes in directivity [66].

Note that GFVV is (ideally) agnostic with regards to  $\hat{s}(f)$  and  $\hat{\chi}(f)$ , hence, we can rewrite (3) as:

$$\hat{\mathbf{v}}(f) = \frac{\hat{\mathbf{h}}(f)}{\hat{\mathbf{a}}(f)}, \quad (4)$$

where  $\hat{\mathbf{h}}(f)$  is the numerator of the rightmost part of (3), while  $\hat{\mathbf{a}}(f) = \hat{\mathbf{w}}(f)^H \hat{\mathbf{h}}(f)$ . The channel-wise inverse Fourier transform of GFVV, yields its temporal analogue, *i.e.*, GTVV:

$$\mathbf{v}(t) = \mathcal{F}^{-1}(\hat{\mathbf{v}}(f)) = \mathbf{h}(t) * a^{-1}(t), \quad (5)$$

with  $(a * a^{-1})(t) = \delta(t)$ .

It is noteworthy that the standard RTF in SH domain [32] (*i.e.*, FDVV), for which the reference is the first (omnidirectional) Ambisonic channel, is a special case of GFVV obtained by setting  $\mathbf{w} = [1 \ 0 \ 0 \ \dots \ 0]^T$ . Nonetheless, it would be preferential to use a beamformer steered approximately towards DoA of the source, as discussed later in this section. This also clarifies the notion of “generality” in GFVV - its reference channel does not correspond to one particular HOA channel, but is a linear combination of all available channels. Likewise, TDVV (relative impulse response in SH domain) becomes a special case of GTVV. We remind the reader that the pseudointensity vector [53], [52] corresponds to the real part of RTF for the FOA signals. In [7], we have argued that this classical DoA estimator is biased in the presence of strong reflections, but, without providing a detailed explanation. We take the opportunity to elaborate this claim in Appendix A.

## B. GTVV estimation

The expression (3) cannot be used directly, even in the noiseless setting, since we expect  $\hat{\mathbf{e}}(f)$  to contain the diffuse late reverberation components. Moreover, a practical estimation method needs to be robust to low SNR levels. The following computationally efficient estimator has been proposed in [31], [22], and represents an adaptation of the well-known estimator based on speech signal nonstationarity [67], [55]. From (3), we have that a GFVV entry  $\hat{v}_i(f)$  could be seen as the ratio between “denoised” versions of  $\hat{b}_i(f)$  and the reference:

$$\hat{v}_i(f) = \frac{\hat{b}_i(f) - \hat{e}_i(f)}{\sum_{\nu} \hat{w}_{\nu}^*(f) (\hat{b}_{\nu}(f) - \hat{e}_{\nu}(f))}. \quad (6)$$

Rearranging the terms in the expression above gives

$$\hat{b}_i(f) = \hat{v}_i(f) \sum_{\nu} \hat{w}_{\nu}^*(f) \hat{b}_{\nu}(f) + \hat{n}_i(f) \quad (7)$$

where we denote  $\hat{n}_i(f) = \hat{v}_i(f) \sum_{\nu} \hat{w}_{\nu}^*(f) \hat{e}_{\nu}(f) - \hat{e}_i(f)$ .

Since the two terms on the right hand side are correlated, we will estimate  $\hat{v}_i(f)$  and noise statistics simultaneously, in the least-squares sense, as originally proposed in [67]. First, the signal is analyzed in time-frequency – particularly, Short-time Fourier transform (STFT) – domain:

$$\hat{b}_i(f, t) = \hat{v}_i(f, t) \sum_{\nu} \hat{w}_{\nu}^*(f) \hat{b}_{\nu}(f, t) + \hat{n}_i(f, t), \quad (8)$$

<sup>2</sup>By the virtue of Paley-Wiener theorem [65].

where  $(f, t)$  designates discrete frequency and time frame indices, respectively, with a slight abuse of notation. Multiplying both sides by  $\hat{b}_l^*(f, t)$ , and taking expectation yields

$$\hat{\phi}_l^2(f, t) = \hat{v}_l(f, t) \sum_{l'} \hat{w}_l^*(f) \hat{\phi}_{l', l}(f, t) + \hat{\sigma}_l(f, t), \quad (9)$$

where  $\hat{\phi}_l^2(f, t) = \mathbb{E}[|\hat{b}_l(f, t)|^2]$  is the variance of the  $l^{\text{th}}$  channel, while  $\hat{\phi}_{l', l}(f, t) = \mathbb{E}[\hat{b}_{l'}(f, t) \hat{b}_l^*(f, t)]$  and  $\hat{\sigma}_l(f, t)$  are the cross-correlations between channels  $l'$  and  $l$ , and between the noise term  $\hat{\eta}_l$  and  $\hat{b}_l$ , respectively. Under the assumption that the noise statistics  $\hat{\sigma}_l$  and the acoustics of the environment (thus, GFVV) evolve slower than speech statistics, and by rearranging the terms, the expression (9) is approximated by

$$\hat{\phi}_l^2(f, t) \approx \begin{bmatrix} \hat{\phi}_{:, l}(f, t)^H \hat{\mathbf{w}}(f) & \mathbf{1} \end{bmatrix} \begin{bmatrix} \hat{v}_l(f) \\ \hat{\sigma}_l(f) \end{bmatrix}, \quad \text{where} \quad (10)$$

$$\hat{\phi}_{:, l}(f, t) = [\hat{\phi}_{0, l}(f, t) \quad \hat{\phi}_{1, l}(f, t) \quad \dots \quad \hat{\phi}_{(L+1)^2-1, l}(f, t)]^T. \quad (11)$$

The approximation above is assumed to hold for a collection of frames centered at  $t_0$ , *i.e.* within  $t \in [t_0 - T/2, t_0 + T/2]$ , for  $T \in 2\mathbb{N}$ . Note that  $\hat{v}_l(f)$  and  $\hat{\sigma}_l(f)$  are assumed constant for this set of frames, which is compactly written as

$$\hat{\phi}_l^2(f) \approx \begin{bmatrix} \hat{\Phi}_{:, l}(f)^H \hat{\mathbf{w}}(f) & \mathbf{1} \end{bmatrix} \begin{bmatrix} \hat{v}_l(f) \\ \hat{\sigma}_l(f) \end{bmatrix}, \quad \text{where} \quad (12)$$

$$\hat{\phi}_l^2(f) = [\hat{\phi}_l^2(f, t_0 + T/2) \quad \dots \quad \hat{\phi}_l^2(f, t_0 - T/2)]^T, \quad (13)$$

$$\hat{\Phi}_{:, l}(f) = [\hat{\phi}_{:, l}(f, t_0 + T/2) \quad \dots \quad \hat{\phi}_{:, l}(f, t_0 - T/2)] \quad (14)$$

and  $\mathbf{1}$  is the all-one vector. This is an overdetermined linear system that can be solved efficiently in the least squares sense (amounts to the inversion of a  $2 \times 2$  matrix), providing an estimate of  $\hat{v}_l(f)$  at frame  $t_0$ .

The quality of the GFVV estimate depends on a number of factors, including the STFT parameters (window type and length, overlap percentage), neighborhood size  $T$ , spectral contents of the excitation signal and noise, and, naturally, room acoustics. We will see later that the dominant wavefront is positioned at the zero-delay index of the RdRIR representation, *i.e.*, in the middle of the time series. Hence, when the dominant wavefront is due to direct sound, capturing the early echoes requires the STFT frame length to be at least twice the mixing time value  $\bar{\tau}_{\max}$ . The estimator presents certain advantages and drawbacks. Conveniently, it requires only the information about the activity of the target sound source. This is also related to its susceptibility to directional interference, thus, a reliable voice activity detector (VAD) is implied. While it can adapt to changes in the acoustic environment, its performance tend to degrade in highly dynamic scenarios. On the other hand, if more refined information is available, such as the noise covariance matrix, one could conceive adaptations of other well-known RTF estimation techniques, *e.g.*, the covariance subtraction and covariance whitening methods [68], [32].

In addition to the previously discussed uncertainties, another unknown parameter is the DoA of the source of interest, which is often a required parameter to design the beamformer vector  $\hat{\mathbf{w}}(f)$ . The following subsection is dedicated to the derivation of a closed-form expression of GTVV, where we demonstrate that – under certain conditions – GTVV can be used to

---

### Algorithm 1 Self-steering GTVV estimator at a frame $t_0$

---

**Require:** STFT tensor  $\{\hat{\mathbf{b}}(f, t) \mid t \in t_0 + [-T/2, T/2], f \in [0, K), l \in [0, (L+1)^2)\}$ , parametric SH dictionary  $\{\mathbf{y}(\theta, \phi)\}_{(\theta, \phi)}$ , num\_iter

**Compute:**  $\hat{\phi}_l^2(f, t) = |\hat{b}_l(f, t)|^2$  and  $\hat{\phi}_{l', l}(f, t) = \hat{b}_{l'}(f, t) \hat{b}_l^*(f, t)$

**Assemble:**  $\hat{\phi}_l^2(f)$  and  $\hat{\Phi}_{:, l}(f)$  from eq. (13) and (14)

Set  $\mathbf{w} = [1 \ 0 \ 0 \ \dots \ 0]^T$

**for** iter  $\in [1, \text{num\_iter}]$  **do**

**for**  $l \in [0, (L+1)^2 - 1]$  **do**

$\hat{v}(f) \leftarrow$  solve (12) for each  $f$

$\mathbf{v}(\tilde{t}) \leftarrow \mathcal{F}^{-1}(\hat{v}(f))$

**end for**

$(\theta_0, \phi_0) \leftarrow \operatorname{argmax}_{(\theta, \phi)} \mathbf{v}(\tilde{t} = 0)^T \mathbf{y}(\theta, \phi)$

$\mathbf{w} \leftarrow \mathbf{y}(\theta_0, \phi_0) / (L+1)^2$

**end for**

**Return:**  $\mathbf{v}(\tilde{t}), \hat{v}(f), (\theta_0, \phi_0)$

---

directly infer the source's DoA. However, we have observed that even if these conditions do not hold, the GTVV vector  $\mathbf{v}(t = 0)$  is usually a good approximation of the SH encoding vector parametrized by a steering angle pointed in the vicinity of DoA. We exploit this observation to devise a heuristic scheme that improves the DoA estimate iteratively. Indeed, a “well-behaved” GTVV representation maintains  $\mathbf{v}(t = 0)$  that is invariant to the slight changes in  $\hat{\mathbf{w}}(f)$ , *i.e.*, it should always point towards DoA. We, therefore, use the current DoA estimate to (re-)steer the beamformer and compute a new GTVV representation. Starting with the omnidirectional reference (the classical relative IR), we monitor the difference in DoA between iterations to deduce whether GTVV has “converged”. Typically, this procedure requires no more than ten iterations. For reader's convenience, its pseudocode is given in Alg. 1 (the specification of the applied signal-independent beamformer is given in Section V).

### C. Closed-form expression

In order to derive a closed-form expression of GTVV, we will treat the numerator and denominator of the GFVV expressions (4) separately.

The numerator  $\hat{\mathbf{h}}(f)$  is simply the “early” part of ATF, normalized by the amplitude of direct component  $\hat{a}_0(f) e^{-j2\pi f \bar{\tau}_0}$ . Hence, its time domain counterpart is the early part of RIR, normalized and shifted to temporal origin - dubbed hereafter *Reduced Room Impulse Response (RdRIR)*:

$$\mathbf{h}(t) = \delta(t) \mathbf{y}_0 + \sum_{n=1}^{N-1} g_n(t - \tau_n) \mathbf{y}_n. \quad (15)$$

Note that, under the assumptions that have been stated earlier,  $\hat{g}_n(f)$  is a real-valued, nonnegative and even function of frequency. Hence, it is easy to show that  $g_n(t) = \mathcal{F}^{-1}(\hat{g}_n(f))$  is also real and even, and that it attains global maximum at  $t = 0$ . Since we also assumed that  $g_n(t)$  has compact support, temporally well-separated wavefronts (having sufficiently dis-

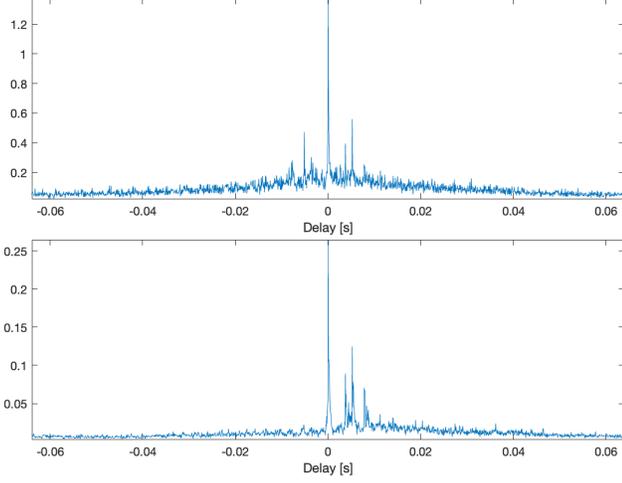


Figure 1: The GTVV estimate's delay-magnitude  $\zeta_v(t)$  without (top) and with (bottom) condition (17) satisfied. The latter is approximately causal, as predicted.

tinct delays  $\tau_n$ ), could be identified by observing peaks of its delay-magnitude time series  $\zeta_h(t)$ :

$$\zeta_h(t) = \|\mathbf{h}(t)\|_2. \quad (16)$$

An explicit solution of  $a^{-1}(t)$  requires more attention. We introduce an additional hypothesis:

$$\left| \sum_{n=1}^{N-1} \hat{\kappa}_n(f) e^{-j2\pi f \tau_n} \right| < 1, \quad \forall f, \quad (17)$$

which is also a sufficient condition for assuring that the impulse response of the reference is a minimum-phase filter [69]. The importance of this condition for the extraction of wavefront parameters will be discussed in the remainder of this section. For the time being, we motivate its interest by visually inspecting two instances of the time series  $\zeta_v(t) = \|\mathbf{v}(t)\|_2$  in Fig. 1, with and without condition (17) satisfied by the vector-valued GTVV function  $\mathbf{v}(t)$ .

The requirement (17) is obviously granted if the magnitude of the direct component is larger than the cumulative magnitude of all reflections in the reference signal:

$$\sum_{n=1}^{N-1} |\hat{\kappa}_n(f)| = \sum_{n=1}^{N-1} \hat{g}_n(f) \left| \hat{\beta}_n(f) \right| < 1, \quad \forall f. \quad (18)$$

Clearly, this is highly unlikely if the reference is the omnidirectional channel. Instead, by using a beamformer steered towards DoA, this hypothesis becomes more and more plausible with the increase in Ambisonic order. Indeed, for beamformers such as maximum-directivity or Minimum Variance Distortionless Response (MVD) [32], having  $|\hat{\mathbf{w}}(f)^H \mathbf{y}_n| < \beta_0$  leads to  $\lim_{L \rightarrow \infty} |\hat{\beta}_n(f)| = 0$ , due to the completeness property of spherical harmonics [66]. Intuitively, with the increase in  $L$ , the spatial response of the beamformer approaches the delta function centered around DoA [32]. Alternatively, one may consider forcing spatial nulls in the directions of strong reflectors (known a priori), *e.g.* by using the Linearly Constrained Minimum Variance (LCMV) [48] beamformer.

Should the condition (17) hold, we can reformulate the GFVV denominator  $\hat{a}^{-1}(f)$  in (4) through the Taylor (geometric) series expansion. In the following, we omit the frequency variable  $f$  for brevity, and let  $\gamma_n := -\hat{\kappa}_n(f) e^{-j2\pi f \tau_n}$ . Then, the denominator in (3) becomes

$$\frac{1}{\hat{a}(f)} = \frac{1}{1 - \sum_{n=1}^{N-1} \gamma_n} = \sum_{k=0}^{\infty} \left( -\sum_{n=1}^{N-1} \gamma_n \right)^k := \sum_{k=0}^{\infty} \lambda_k, \quad (19)$$

where each element  $\lambda_k$  of the last sum is developed using multinomial theorem into

$$\lambda_k = \sum_{i_1+i_2+\dots+i_{N-1}=k} \frac{k!}{i_1!i_2!\dots i_{N-1}!} \prod_{q=1}^{N-1} \gamma_q^{i_q}.$$

Evaluating  $\lambda_k$  for  $k = 0, 1, 2, 3 \dots$  yields:

$$\begin{aligned} \lambda_0 &= 1, \\ \lambda_1 &= \sum_i \gamma_i, \\ \lambda_2 &= \sum_i \gamma_i^2 + 2 \sum_i \sum_{m \neq i} \gamma_i \gamma_m, \\ \lambda_3 &= \sum_i \gamma_i^3 + 3 \sum_i \sum_{m \neq i} \gamma_i^2 \gamma_m + 6 \sum_i \sum_{m \neq i} \sum_{p \neq \{i, m\}} \gamma_i \gamma_m \gamma_p \\ &\dots \end{aligned} \quad (20)$$

where all sums correspond to indices in  $[1, N]$ . Due to directional beamforming, we expect only a subset of reflections to have non-negligible magnitudes  $|\kappa_n(f)| \gg 0$ ,  $n \in [1, N-1]$  (with the size of this subset decreasing with the increase in Ambisonic order, as discussed before). Therefore, the magnitudes of ‘‘cross-terms’’ in the expressions above are more likely to diminish than the leftmost terms that correspond to isolated reflections. We simplify the expression by aggregating all cross terms in a single variable  $\eta(f)$ , hence the GFVV denominator is represented as

$$\sum_{k=0}^{\infty} \lambda_k = 1 + \sum_{i=1}^{N-1} \sum_{k=1}^{\infty} (-\hat{\kappa}_i(f))^k e^{-j2\pi f k \tau_i} + \hat{\eta}(f). \quad (21)$$

In time domain, this yields the following expression:

$$a^{-1}(t) = \delta(t) + \sum_{i=1}^{N-1} \sum_{k=1}^{\infty} (-1)^k \kappa_i^{*k}(t - k\tau_i) + \eta(t), \quad (22)$$

where  $\kappa_i^{*k}(t) = \mathcal{F}^{-1}(\hat{\kappa}_i(f)^k) = \overbrace{(\kappa_i * \kappa_i * \dots * \kappa_i)}^{k-1 \text{ convolutions}}(t)$ .

Plugging (15) and (22) into (5), and manipulating the terms within, produces

$$\mathbf{v}(t) = \delta(t) \mathbf{y}_0 + \sum_{n=1}^{N-1} \sum_{k=1}^{\infty} (-1)^k \kappa_n^{*k}(t - k\tau_n) * (\mathbf{y}_0 \delta(t) - \mathbf{y}_n \beta_n^{-1}(t)) + \tilde{\eta}(t), \quad (23)$$

where  $\tilde{\eta}(t)$  again accounts for  $\eta(t)$ , augmented by additional cross-convolutions between different reflections. One can make several observations of the GTVV representation

(23). First, due to the assumed compact support of  $\kappa_n(t)$ , GTVV is approximately causal, *i.e.*,  $v(t < 0) \approx \mathbf{0}$ . Second, if the stronger condition (18) holds, we expect GTVV to be somewhat sparse (as the energy of  $\kappa_n^{**k}(t)$  decreases with  $k$ ). Third, interestingly, GTVV still allows us to immediately identify the SH vector corresponding to direct component, by evaluating  $v(t = 0)$ , as for RdRIR in (15). However, even by neglecting the cross-terms  $\tilde{\eta}(t)$ , it is obvious that the remainder of the GTVV expression is more complex, presenting itself as a series of repeated convolutions with alternating sign, for each wavefront  $n$ .

Hence, without additional assumptions, we cannot easily identify the remaining wavefronts. For instance, if we conjecture that the initial terms ( $k = 1$ ) of each series are not strongly affected by another series *and* that  $\beta_n(0) \ll 1$ , then the largest peaks of  $\zeta_v(t)$  would likely<sup>3</sup> correspond to SH vectors  $\mathbf{y}_n$ , *i.e.*

$$v(\tau_n) = g_n(0) (\mathbf{y}_n - \beta_n(0)\mathbf{y}_0) \approx g_n(0)\mathbf{y}_n. \quad (24)$$

If  $\beta_n(0) \ll 1$  does not hold, yet we still assume that  $v(\tau_n)$  could be isolated (*e.g.*, for the strong reflections), we can exploit the fact that  $\mathbf{y}_0$  can be pre-estimated from  $v(t = 0)$  to estimate the wavefront vector  $\mathbf{y}_n$ . In this case, we consider only a wideband beamformer  $\mathbf{w}$  and propose to solve a nonlinear optimization problem:

$$(\theta_n, \phi_n) = \underset{(\theta, \phi)}{\operatorname{argmin}} v(\tau_n)^\top (\mathbf{I} - \mathbf{y}_0 \mathbf{w}^\top) \mathbf{y}(\theta, \phi), \quad (25)$$

where  $\mathbf{y}(\theta, \phi)$  is a SH vector for the given azimuth and elevation parameters. In the cost function (25), we have used the expression for the (constant) spatial response of a wideband beamformer  $\beta_n = \mathbf{w}^\top \mathbf{y}_n$ . To avoid explicitly solving the optimization problem above, one may use a dictionary of normalized SH encoding vectors  $\mathbf{y}(\theta, \phi)$ , parametrized from a discrete grid of directions  $\{(\theta, \phi)\}$ , and choose the atom most correlated with  $v(\tau_n)^\top (\mathbf{I} - \mathbf{y}_0 \mathbf{w}^\top)$ . Unfortunately, the matrix  $\mathbf{I} - \mathbf{y}_0 \mathbf{w}^\top$  is not invertible (otherwise, one could directly obtain an estimate of  $\mathbf{y}_n$ ), which is easy to show by applying the Sherman-Morrison formula [70].

Finally, we remark that the presented derivation does not directly depend on the DoA direction  $\mathbf{y}_0$  used for the distortionless constraint  $\mathbf{w}(f)^\mathbf{H} \mathbf{y} = \text{const}$ , as long as the Taylor series condition (17) holds. In other words, if a sufficiently selective beamformer is focused on some other wavefront  $\mathbf{r}$  (*e.g.*, a dominant reflection), one would replace  $\mathbf{y}_0$  by  $\mathbf{y}_r$ , and a similar expression applies:

$$v(t) = \delta(t)\mathbf{y}_r + \sum_{n \neq r} \sum_{k=1}^{N-1} \sum_{k=1}^{\infty} (-1)^k \kappa_n^{**k}(t - k\tau_n) * (\mathbf{y}_r \delta(t) - \mathbf{y}_n \beta_n^{-1}(t)) + \tilde{\eta}(t), \quad (26)$$

except that the quantities  $g_n$  and  $\tau_n$  are now relative to the absolute gain  $a_r(f)$  and ToA  $\bar{\tau}_r$  of this wavefront, respectively. As a consequence, relative gains  $g_n$  would not be bounded by 1, relative delays  $\tau_n = \bar{\tau}_n - \bar{\tau}_r$  could have negative values, and  $v(t = 0)$  would encode the reflection direction  $(\theta_r, \phi_r)$ .

<sup>3</sup>It may still happen that later terms ( $k > 1$ ) of dominant reflections have larger peaks than the initial terms of weaker wavefronts!

Compared to the GTVV computed using DoA, this variant would be “shifted” to the left by  $|\tau_0|$ .

#### IV. ESTIMATION OF REDUCED RIR

We have argued that GTVV is better adapted to reverberant acoustic conditions than the “standard” relative impulse response for which the reference signal is the zero-order Ambisonic channel. Nevertheless, it is still limited by the spatial selectivity of the applied beamformer - for example, if the signal-independent maximum directivity beamformer is used, its directivity will be proportional to the square of Ambisonic order [32]. However, affordable Ambisonic microphone arrays usually do not provide very high order Ambisonic formats - most often, they are only capable of recording the FOA signals [34]. Furthermore, the frequency support of higher order channels progressively decreases with the HOA order, as noise amplification at low frequencies, and spatial aliasing at high frequencies start to kick-in [71]. Unfortunately, the favorable theoretical properties of GTVV tend to diminish at low Ambisonic orders, due to the inability of the applied beamformer to effectively suppress the reflections. The problem is further exacerbated with the increase in the microphone-to-source distance, since more reflections fall within the main lobe of the beamformer. In practice, we observe that the GTVV imprint is no longer causal (as seen in Fig. 1), and that the estimated directions are less accurate.

Moreover, even when the GTVV expression (23) remains valid, identifying the directions and delays by peak-picking is not straightforward, as discussed in the previous section. In fact, such a “well-behaved” GTVV can be seen as the reduced RIR (15), convolved by the minimum-phase filter (22). The consequence is that the same reflection is infinitely “echoed” at the time instances corresponding to integer multiples of its relative delay, with the alternating sign and the decreasing magnitude. Thus, these series can interfere with one another, altering the information within, or even masking the presence of weaker reflections. Undoubtedly, it is much easier and intuitive to extract information directly from RdRIR (15). In this section, we propose a simple method to estimate the latter from the observed GTVV time series, even if the convergence condition (17) is not satisfied. The development is based on the celebrated Padé-Prony method for the pole-zero modeling [72], which is very similar to traditional Autoregressive Moving Average (ARMA) model for stochastic time series.

In the following, we consider a beamformer steered towards DoA, since the same method could be straightforwardly adapted when other wavefronts are considered. We start by rewriting (5) as

$$(\mathbf{v} * a)(t) = \mathbf{h}(t), \quad (27)$$

and recall that  $a(t) = \mathcal{F}^{-1}(\mathbf{w}(f)^\mathbf{H} \mathbf{h}(f))$ , *i.e.*,

$$\begin{aligned} a(t) &= \mathcal{F}^{-1} \left( 1 + \sum_{n=1}^{N-1} \hat{\kappa}_n(f) e^{-j2\pi f \tau_n} \right) \\ &= \delta(t) + \sum_{n=1}^{N-1} \kappa_n(t - \tau_n). \end{aligned} \quad (28)$$

Since we have assumed that all  $\kappa_n(t)$  have compact support,  $a(t)$  is a causal filter (but, not necessarily a minimum-phase!). We already know from (15) that RdRIR  $\mathbf{h}(t)$  is a causal vector sequence, *i.e.*,  $\mathbf{h}(t < 0) \approx \mathbf{0}$ . Moreover, RdRIR has finite support - beyond the relative delay  $\tau_{\max} = \bar{\tau}_{\max} - \bar{\tau}_0$ , *i.e.*, relative to the mixing time  $\bar{\tau}_{\max}$ , we expect  $\mathbf{h}(t > \tau_{\max}) \approx \mathbf{0}$  to hold as well (being the feature of the “denoising” estimator presented in subsection III-B). Likewise, one may argue that the early part of RIR (hence, RdRIR) is relatively sparse [48], [43], [44]. The filter  $a(t)$  would be even sparser, as we expect the beamforming operation to suppress certain reflections in the reference signal. Finally, one may observe from (15) that, for any  $t$ , the zero-order entry of  $\mathbf{h}(t)$  is non-negative. Our aim is to take advantage of all this prior knowledge to estimate  $a(t)$  directly from  $\mathbf{v}(t)$ , and then extract  $\mathbf{h}(t)$  by convolving its estimate with GTVV, as in (27).

As usual in digital signal processing, we do not handle continuous functions  $\mathbf{v}(t)$ ,  $\mathbf{h}(t)$  and  $a(t)$ , but their sampled versions. We make a leap of faith and assume that the latter are not substantially affected by aliasing, meaning that the properties discussed above are generally preserved. We denote by  $j$  the time sample index taking values<sup>4</sup> in  $[-J/2 + 1, J/2]$ , within an STFT frame of length  $J \in 2\mathbb{N}$ . Thus, both GTVV and RdRIR are represented by the real-valued matrices  $\mathbf{V}$  and  $\mathbf{H}$  of size  $(L + 1)^2 \times J$ , *i.e.* their columns  $\mathbf{v}_{:,j}$  (accordingly,  $\mathbf{h}_{:,j}$ ) are akin to evaluating  $\mathbf{v}(t)$  and  $\mathbf{h}(t)$  at some time instant  $t$ . Analogously, the rows  $\mathbf{v}_{l,:}$  (accordingly,  $\mathbf{h}_{l,:}$ ) correspond to the  $l^{\text{th}}$  channels of the two representations. The filter  $a(t)$  is replaced by a vector  $\mathbf{a} \in \mathbb{R}^{j_{\max} + 1}$ , where the hyperparameter  $j_{\max}$  denotes the sample index corresponding to  $\tau_{\max}$ , *i.e.*, the assumed relative delay of the “last” wavefront in RdRIR. The coefficients  $a_{j < 0}$  and  $a_{j > j_{\max}}$  are assumed to be zero, hence, these are not included in the estimation vector  $\mathbf{a}$ .

Setting aside the sparsity hypothesis for now, note that enforcing  $\zeta_{\mathbf{h}}(t) = 0$ , for  $t < 0$  and  $t > \tau_{\max}$ , amounts to minimizing the following cost function:

$$\min_{\mathbf{a}} \sum_{l=0}^{(L+1)^2-1} \sum_{j \notin [0, j_{\max}]} (v_{l,:} * \mathbf{a})_j^2, \quad \text{s.t. } a_0 = 1. \quad (29)$$

The equality constraint is due to  $a(0) = \delta(0)$  in (28), with the Dirac delta distribution replaced by the Kronecker delta function in the discrete version. We remark that (29) is a particular multichannel linear prediction problem, with the filter  $\mathbf{a}$  being common for all channels  $l \in [0, (L + 1)^2 - 1]$ . This is advantageous - since the problem is overdetermined, the estimate of  $\mathbf{a}$  should be more resilient to GTVV estimation errors and noise. Furthermore, the estimation should become more accurate as the channel order  $L$  increases.

There are multiple approaches of addressing linear prediction problems, but probably the most well-known are the autocorrelation method and the covariance method [72]. In both cases, solving the constrained quadratic problem comes

down to a linear system, compactly written as

$$\sum_{j=1}^{j_{\max}} a_j r(j, \mathbf{s}) = -r(0, \mathbf{s}), \quad (30)$$

$$\text{where } r(j, \mathbf{s}) = \sum_{l=0}^{(L+1)^2-1} r_l(j, \mathbf{s}) = \sum_{l=0}^{(L+1)^2-1} \sum_{j'} v_{l,j'-j} v_{l,j'-\mathbf{s}}. \quad (31)$$

The two methods differ in the way they deal with the signal edges, *i.e.*, how they define the range of the summation variable  $j'$ . The autocorrelation method applies zero-padding ( $v_{l,j'} = 0$ , for  $j' \in [0, j_{\max}]$  and  $j' \notin [-J/2 + 1, J/2]$ ), while the covariance method considers only valid parts of the convolution (where the two sequences overlap and  $j' \notin [0, j_{\max}]$ ), and discards the rest. Therefore, the coefficients  $r(j, \mathbf{s})$  would be somewhat different, yielding different solutions. Particularly, the filter  $\mathbf{a}$  estimated by the autocorrelation method is always minimum-phase [72], but the corresponding linear system has Toeplitz structure, hence it can be solved by the Levinson-Durbin algorithm [70] with  $O((j_{\max} + 1)^2)$  time complexity. This is significantly more efficient compared to  $O((j_{\max} + 1)^3)$  of the covariance method.

Furthermore, calculating the coefficients (31) of the normal equations (30) generally requires  $O((L + 1)^2 J^2)$  multiplications, but for the autocorrelation method this cost is reduced, thanks to the duality of autocorrelation and power spectrum [72]. Due to the symmetry property of autocorrelation, we have  $r_l(j, \mathbf{s}) = r_l(j - \mathbf{s}) = r_l(\mathbf{s} - j)$ . Now define

$$v_{l,j}^- = \begin{cases} v_{l,j}, & j \in [-J/2 - 1, 0), \\ 0, & \text{otherwise,} \end{cases} \quad (32)$$

$$v_{l,j}^+ = \begin{cases} v_{l,j}, & j \in (j_{\max}, J/2], \\ 0, & \text{otherwise,} \end{cases} \quad (33)$$

and let

$$r_l^-(j - \mathbf{s}) = \sum_{j'=-J/2-1}^{J/2} v_{l,j'-j}^- v_{l,j'-\mathbf{s}}^- = \mathcal{F}^{-1} (|\hat{v}_{l,:}^-|^2)_{j-\mathbf{s}}, \quad (34)$$

$$r_l^+(j - \mathbf{s}) = \sum_{j'=-J/2-1}^{J/2} v_{l,j'-j}^+ v_{l,j'-\mathbf{s}}^+ = \mathcal{F}^{-1} (|\hat{v}_{l,:}^+|^2)_{j-\mathbf{s}}, \quad (35)$$

where  $\hat{v}_{l,:}^-$  and  $\hat{v}_{l,:}^+$  are the frequency representations of  $v_{l,:}^-$  and  $v_{l,:}^+$ , respectively. Having  $r_l(j - \mathbf{s}) = r_l^-(j - \mathbf{s}) + r_l^+(j - \mathbf{s})$  and

$$\begin{aligned} r(j, \mathbf{s}) &= \sum_{l=0}^{(L+1)^2-1} r_l(j - \mathbf{s}) \\ &= \mathcal{F}^{-1} \left( \sum_{l=0}^{(L+1)^2-1} (|\hat{v}_{l,:}^-|^2 + |\hat{v}_{l,:}^+|^2) \right)_{j-\mathbf{s}}, \end{aligned} \quad (36)$$

we obtain the multichannel linear prediction coefficients at  $O((L + 1)^2 J \log J)$  computational cost.

We note that the presented approach is a variant of classical Prony-like estimation, lauded for its computational efficiency, yet a more elaborate technique may be applied. For instance,

<sup>4</sup>We intentionally permit negative indexing, to preserve the intuition that the temporal dimension is centered at zero.

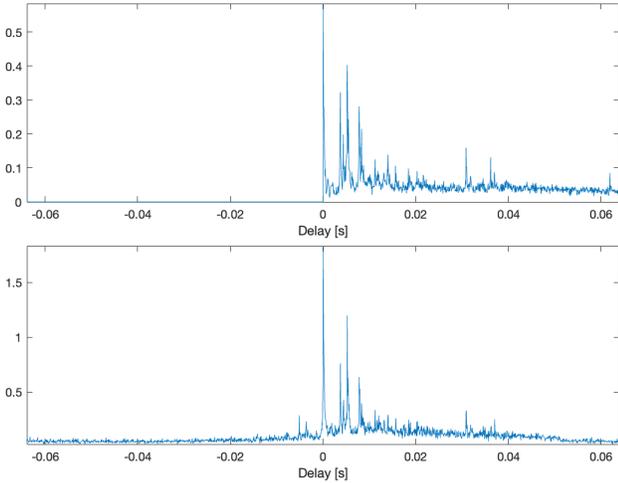


Figure 2: Delay-magnitude  $\zeta_{\hat{h}}(t)$  of the ground truth (top) and estimated (bottom) RdRIR, recovered from the acausal GTVV representation given in Fig 1 (top).

one could perform alternating minimization to improve the estimates of  $a(t)$  and  $\mathbf{h}(t)$  iteratively, in the spirit of the Steiglitz-McBride algorithm [73]. Therefore, to incorporate the sparsity and non-negativity assumptions, we propose to jointly optimize the two variables:

$$\begin{aligned} \min_{\mathbf{a}, \mathbf{H}} \sum_j \|\mathbf{h}_{:,j}\|_2 \\ \text{s.t. } \mathbf{h}_{l,:} = \mathbf{v}_{l,:} * \mathbf{a}, \\ h_{\forall l,j < 0} = h_{\forall l,j > j_{\max}} = 0, \\ h_{0,\forall j} \geq 0 \text{ and } a_0 = 1. \end{aligned} \quad (37)$$

However, the imposed modeling constraints make this problem inconsistent in practice (*e.g.* due to estimation errors and noise). While one could reformulate the problem such that the constraints are relaxed - for instance, by introducing a squared norm penalty instead of the first equality constraint, it would require introducing a new regularization hyperparameter. Instead, we propose using Alternating Directions Method of Multipliers (ADMM), a first-order optimization framework based on Douglas-Rachford splitting [74]. ADMM is particularly effective for optimization problems involving linearly dependent variables, provided that the solutions of intermediate optimization problems are efficiently obtained. Another convenient feature of ADMM is that, in the inconsistent setting, its iterates could produce the best approximation pair, *i.e.*, a pair of estimates of  $\mathbf{a}$  and  $\mathbf{H}$  for which the residuals  $\mathbf{h}_{l,:} - \mathbf{v}_{l,:} * \mathbf{a}$  attain the lowest norm [75]. In Appendix B, we instantiate ADMM for the problem (37) - an interested reader can easily derive the algorithm by following the tutorial article [76] by Boyd et al. Fig. 2 illustrates an example of the reconstructed RdRIR using the proposed ADMM.

## V. EXPERIMENTS

We evaluate RdRIR estimation on data generated using simulated and recorded Ambisonic RIRs. In particular, the

autocorrelation (AC), covariance (COV) and ADMM methods applied to GTVV are benchmarked. As baselines, we use the canonical MCLMS algorithm, as well as the noise-robust multichannel frequency-domain least mean squares (RNMCFLMS) [77] version. Both least mean squares (LMS) implementations are from the *Blind System Identification and Equalization (BSIE)* toolbox [78]. As additional baselines, we use the “plain” GTVV and TDVV representations.

The “fully blind” scenario is assured in all experiments, *i.e.*, the algorithms only have access to the (noisy) observed signals, which are generated by convolving the multichannel RIRs with 10 s of speech data, taken from the publicly available LibriSpeech corpus [79]. The TDVV and GTVV representations are estimated from these measurements using the approach described in III-B, with the latter being obtained through the “self-steering” heuristics, explained at the end of the same subsection. The number of frames  $T$  used for the estimation corresponds to the 0.5 s buffer. The GTVV reference signal is obtained using a “regular” (or “Plane-Wave Decomposition” or “Maximum-Directivity”) beamformer [32] pointing towards the iteratively estimated main DoA  $(\hat{\theta}_0, \hat{\phi}_0)$ , *i.e.*:  $\mathbf{w} = \mathbf{y}(\hat{\theta}_0, \hat{\phi}_0) / (L + 1)^2$ , assuming that the spherical harmonic function basis is 3D-Normalized [80]. Three HOA orders are considered:  $L \in \{1, 2, 3\}$ , while the common sampling rate is set to  $f_s = 16$  kHz. The STFT representation is computed by applying Tukey window of length 0.128 s, with 75% overlap between succeeding frames. The Hendriks algorithm [81] is used as VAD. When applied to the *clean* speech data, this algorithm estimates that about 50% of all STFT frames are voiced.

Since the RdRIR representation is invariant to global scale and ToA offset, it cannot be directly compared to the ground truth multichannel RIRs. Moreover, since the measured RIRs are sampled versions of continuous impulse responses, some of their segments may have undergone sign inversion (due to the action of the anti-aliasing filter). Assuming the same filter has been applied to all channels, it suffices to observe the sign of the omnidirectional component, and then, if the latter is negative, change the signs of all channels at the given sample. Following the sign correction, RIRs are shifted to temporal origin, such that the strongest wavefront - which is assumed to correspond to the direct propagation - is located at  $t = t_0 = 0$ . Finally, the obtained multichannel sequence  $\mathbf{h}(t)$  is rescaled such that the vector at  $t = 0$  is of unit magnitude, and is hereafter referred to as ground truth RdRIR.

The chosen evaluation metrics are aimed to reflect the algorithms’ ability to recover directions and relative delays of early echoes. For that reason, we have avoided the common “normalized project misalignment” (NPM) error [78], which is a point-wise metric that indiscriminately penalizes even small temporal deviations from the ground truth. Instead,  $N = 15$  largest peaks of the ground truth delay-magnitude representation (16) are selected, from which the corresponding delays  $\{t_0, t_1, \dots, t_{N-1}\}$  and encoding vectors  $\{\mathbf{h}(t_0), \mathbf{h}(t_1), \dots, \mathbf{h}(t_{N-1})\}$  are logged. We independently apply the same peak-picking procedure to the delay-magnitude representation  $\zeta_{\hat{h}}(t)$  of a given estimator  $\hat{\mathbf{h}}(t)$ , which yields another set of delays  $\{t_0, \hat{t}_1, \dots, \hat{t}_{N-1}\}$ , associated with en-

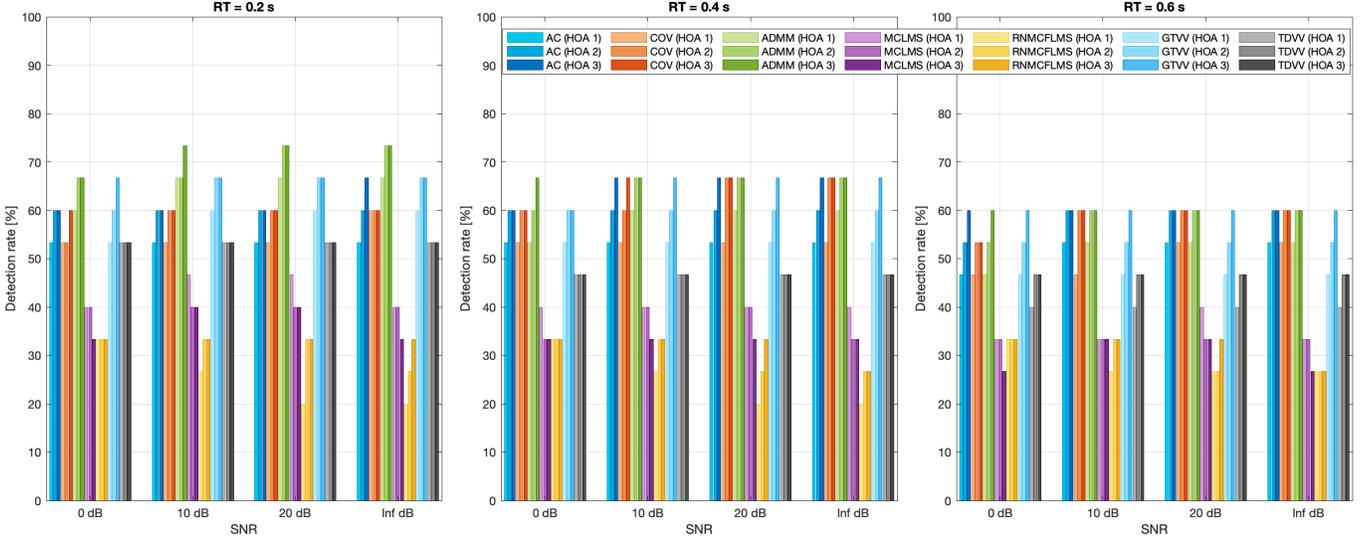


Figure 3: Median detection rate  $\tilde{P}$  results for all methods, with respect to the RT and SNR settings.

coding vectors  $\{\tilde{\mathbf{h}}(t_0), \tilde{\mathbf{h}}(\tilde{t}_1), \dots, \tilde{\mathbf{h}}(\tilde{t}_{N-1})\}$  (note that the zero-delay vector, akin to the DoA delay  $t_0$ , is always retained). We keep only  $\tilde{N}$  of those estimates  $\tilde{\mathbf{h}}_j := \tilde{\mathbf{h}}(\tilde{t}_j)$  whose delays  $\tilde{t}_j$  are within a five-sample temporal neighborhood of the ground truth wavefronts - hence, the relative delay error tolerance is  $|t_i - \tilde{t}_j| \leq 0.3$  ms. The percentage of retained estimates

$$\tilde{P} = \frac{\tilde{N}}{N} 100\%,$$

is to be interpreted as a detection rate indicator. Indeed, since the number of retrieved peaks is equal for both ground truth and the estimate, detection precision and recall have the same value. For all ground truth wavefronts and retained estimates, we then find the closest (in the least squared sense) SH vectors  $\mathbf{y}_i := \mathbf{y}(\theta_i, \phi_i)$ , respectively  $\tilde{\mathbf{y}}_j = \tilde{\mathbf{y}}(\tilde{\theta}_j, \tilde{\phi}_j)$ , parametrized by the appropriate azimuth and elevation values. We use these directions to evaluate angular errors between an estimate and the associated ground truth wavefront:

$$\tilde{\epsilon}_{i,j} = \angle \left( (\theta_i, \phi_i), (\tilde{\theta}_j, \tilde{\phi}_j) \right),$$

where  $\angle(\cdot)$  denotes the great-circle distance for the given pair of directions. Finally, knowing that measured RIRs do not perfectly obey the structure of analytic SH vectors, we also evaluate the coherence  $\tilde{c}_{i,j}$  between a “raw” ground truth vector  $\mathbf{h}_i := \mathbf{h}(t_i)$ , and its estimate  $\tilde{\mathbf{h}}_j$ :

$$\tilde{c}_{i,j} = \frac{\mathbf{h}_i^\top \tilde{\mathbf{h}}_j}{\|\mathbf{h}_i\| \|\tilde{\mathbf{h}}_j\|}.$$

Let us summarize the steps involved:

- 1) Define the STFT frame length to be about  $2\bar{\tau}_{\max}$  and compute the tensor  $\tilde{\mathbf{b}}(f, t)$ .
- 2) Estimate  $\hat{\mathbf{v}}(f)$ ,  $\mathbf{v}(t)$  and DoA through the “self-steering” procedure in Alg. 1 (for TDVV, set num\_iter = 1).
- 3) If RdRIR estimation is used, get  $\tilde{\mathbf{h}}(t)$  from either:
  - AC** - compute coefficients (36) from  $\hat{\mathbf{v}}_{l_i}^-$  and  $\hat{\mathbf{v}}_{l_i}^+$ , assemble and solve the linear Toeplitz system (30).

**COV** - directly compute coefficients (31), assemble and solve the linear system (30).

**ADMM** - set the number of iterations and the parameter  $\mu$ , pre-compute the coefficients (51), iterate the following steps: compute (47) and (52), solve the Toeplitz system (48), compute (44), (45).

If RdRIR is not being estimated, set  $\tilde{\mathbf{h}}(t) = \mathbf{v}(t)$ .

- 4) Compute  $\zeta_{\tilde{\mathbf{h}}}(t) = \|\tilde{\mathbf{h}}(t)\|_2$  and choose the delay indices of  $N$  largest peaks. Preserve only the indices within the prescribed relative delay error tolerance.
- 5) For the retained indices, estimate the direction of each corresponding vector  $\tilde{\mathbf{h}}_j$ :

If RdRIR analysis has been applied, find the parameters  $(\tilde{\theta}_j, \tilde{\phi}_j)$  that minimize the  $\ell_2$  distance of a SH vector-valued function  $\mathbf{y}(\theta, \phi)$  to  $\tilde{\mathbf{h}}_j$ .

If GTVV/TDVV is used directly, find the parameters  $(\tilde{\theta}_j, \tilde{\phi}_j)$  that minimize (25).

- 6) Calculate the performance metrics discussed above.

#### A. Simulated RIRs

In simulated experiments, we mimic the scenario where the microphone array is fixed, while the speech source is mobile, slowly moving at the average speed of 0.97 km/h. This value is only slightly below the average speed of a moving talker in Task 3 of the LOCATA challenge [2]. We have adapted the widely used RIR generator software [82], in order to generate Ambisonic RIRs at each of 100 uniformly sampled positions along 10 randomly generated smooth source trajectories. The corresponding microphone array positions are randomly chosen in the  $xy$ -plane, such that the array altitude is kept fixed at 1.2 m. The microphone signals are obtained by sliding convolution and the spatial interpolation technique implemented in the Roomsimove toolbox [83]. The virtual “room” has dimensions  $5 \times 4 \times 3$  m<sup>3</sup>, while the reverberation time (RT) takes values from  $\{0.2s, 0.4s, 0.6s\}$ . Signals are corrupted by diffuse babble noise whose impulse response has

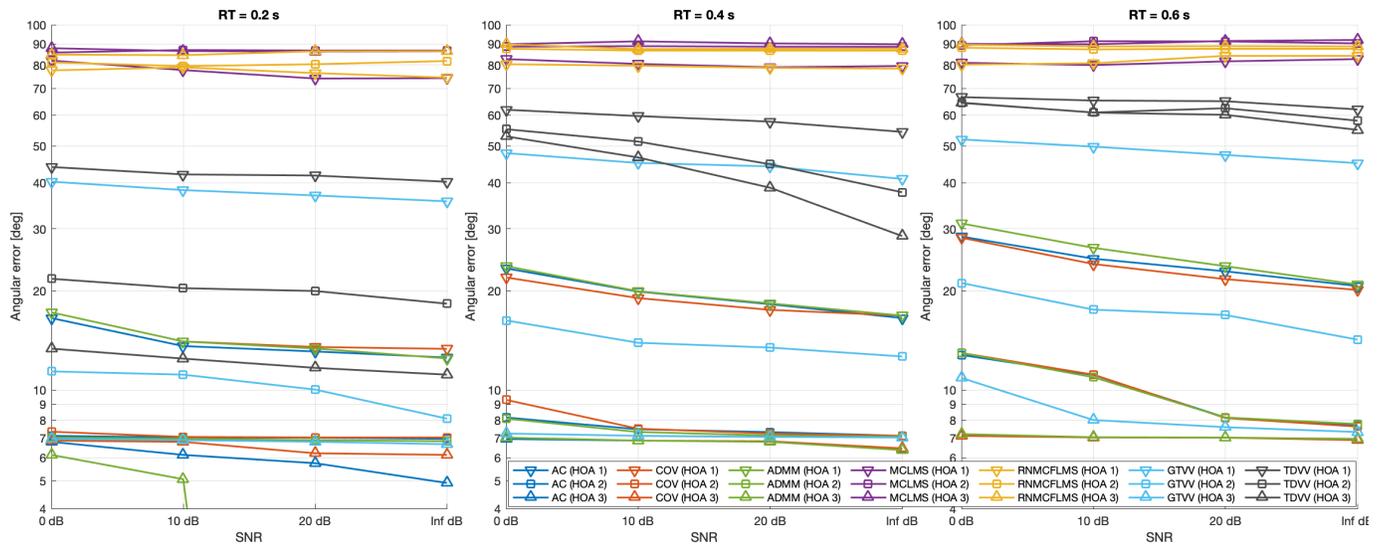


Figure 4: Median angular error for the detected wavefronts, relative to RT and SNR levels. Out-of-scope values (for ADMM (HOA 3) with SNR=20 or Inf dB, and RT=0.2s) are null.

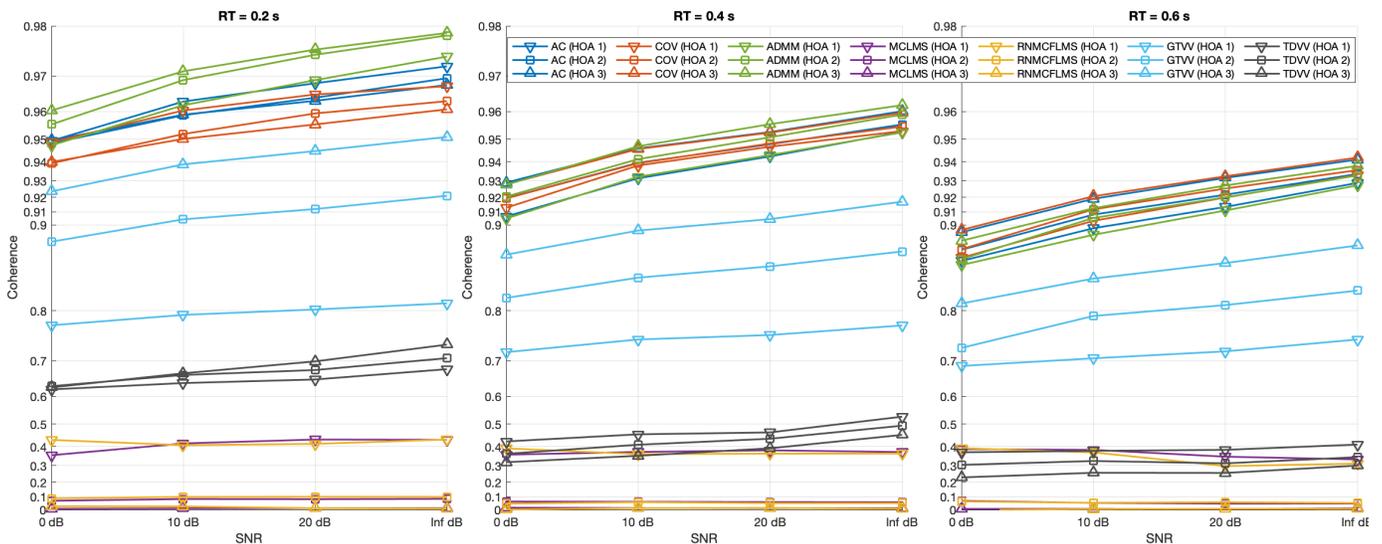


Figure 5: Median coherence for the detected wavefronts, relative to RT and SNR levels.

been obtained by extracting the reverberant parts of several RIRs corresponding to random positions within the room, and then computing their average, as done in [49]. The considered SNR levels are 0, 10, 20 dB and noiseless (“∞” dB setting).

The reported results are computed on the ensemble of generated data, *i.e.* from all generated trajectories. For every performance metric, the experimental outputs are presented as a set of subfigures relative to each RT setting. Within a subfigure, the results are given as a function of the varying SNR level. The detection rate  $\hat{P}$  results are given in the form of bar plots in Fig. 3. The RdRIR-based methods generally provide better detection than the baseline approaches, with the ADMM variant obtaining the highest percentage of accurate detections, most notably at low HOA orders. The angular error and coherence performance on the detected wavefronts are

presented in Fig. 4 and Fig. 5, respectively. One may observe the same trend, with the proposed methods outperforming the baselines, often by a large margin. The LMS baselines performed considerably worse than the other methods, which is also reflected in their poor coherence scores in Fig. 5. They yield essentially similar, quasi-random results across all RT and SNR levels. The obtained results were confirmed by the scores of two sample t-tests [84], evaluated for each pair of estimation methods for a given HOA order. Statistical tests also indicate that, when compared to one another, the proposed RdRIR estimation methods achieve similar performance in terms of the attained angular error and coherence.

With the exception of the LMS baselines, the performance of all tested methods improves with the increase in SNR and HOA order, and drops with the increase in RT. A possible

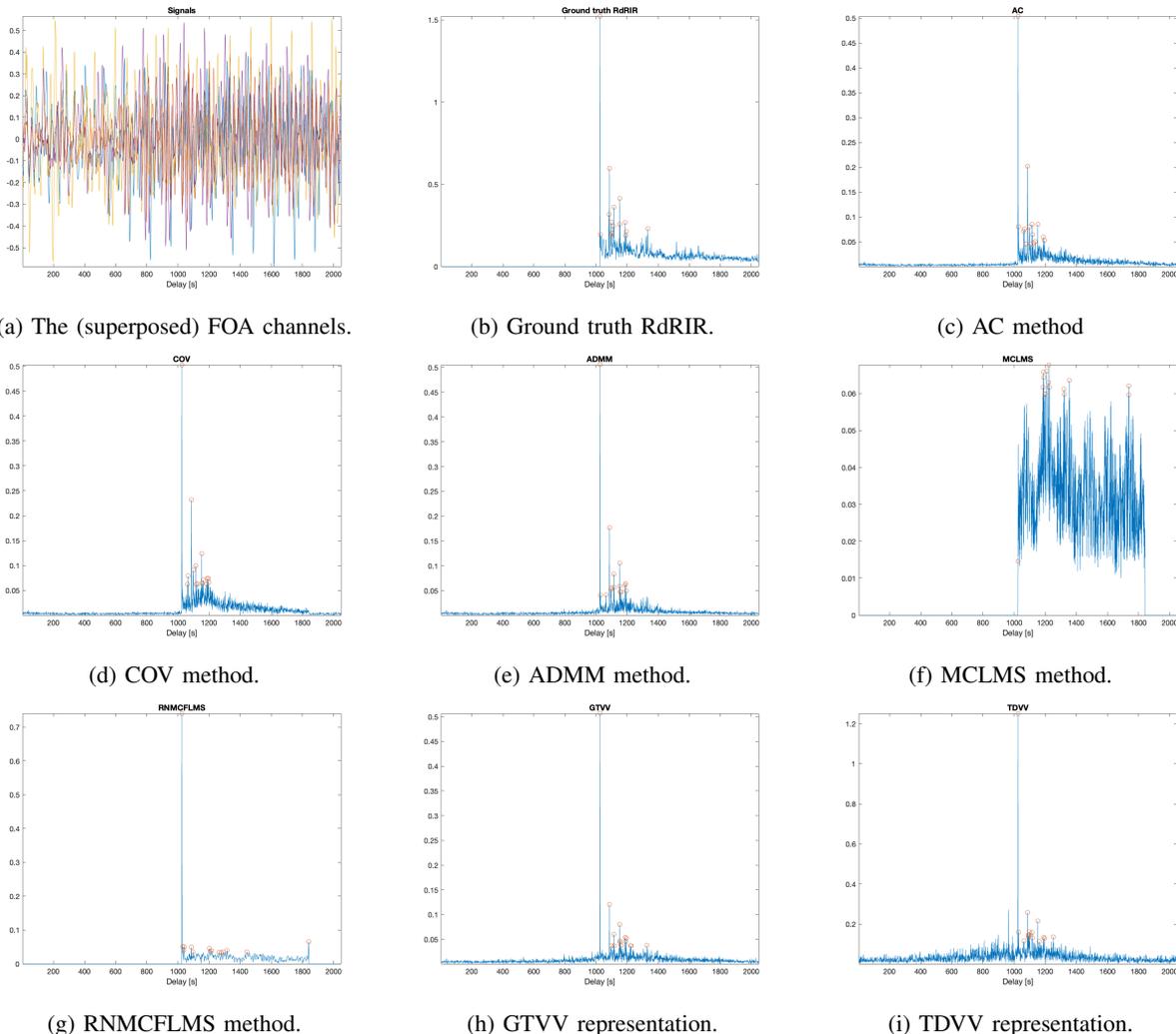


Figure 6: Delay-magnitude representations of the estimates for all tested methods (Figs. 6c-6i), from the FOA recording in Fig. 6a. The ground truth RdRIR is given in Fig. 6b, and the selected peaks are denoted by red circles.

remedy for the latter may be to apply the “channel shortening” technique, *i.e.* to pre-process the input signals by a multichannel dereverberation algorithm (*e.g.* [85]) before the RdRIR estimation.

### B. Recorded RIRs

In order to evaluate the performance of proposed methods in more realistic conditions, we use the dataset of recorded Ambisonic RIRs from University of Aalto, Finland [86]. The dataset contains RIRs in different reverberation conditions, obtained by varying acoustic absorbers in 5 steps, from mild reverberation ( $T_{20}$  around 0.37 s) to highly reverberant ( $T_{20}$  about 1.21 s). The authors have recorded RIRs corresponding to all combinations of 3 positions of a sound source (Genelec 8331A coaxial loudspeaker), and 7 positions of a microphone array (mhAcoustics Eigenmike® em32 and Zylia ZM-1). However, we do not consider the two microphone positions for which the source is facing the opposite direction. Instead, we use these to generate the diffuse impulse response for the additive babble noise, as discussed in the previous

subsection. In this series of experiments, the SNR is fixed to 20 dB (investigating the robustness of the proposed methods under different SNRs, on real data, is left for future work).

For visual comparison, examples of the obtained delay-magnitude representations, along with the corresponding input FOA signals, are presented in Fig. 6 (note that each of the four FOA channels in the subfigure 6a is given in a different color). While for the proposed AC, COV and ADMM methods (and even the GTVV representation), these strongly resemble the ground truth, for TDVV and the LMS baselines this is clearly not the case.

The results are given in Tables II and I, as a function of HOA order,  $T_{20}$  and the type of spherical microphone array (SMA) used for recording RIRs. The best results are emphasized by the boldface font. The proposed methods clearly outperform the baselines, in terms of all evaluation metrics. The three RdRIR-based approaches obtain comparable results, without a clear winner in terms of estimation performance. However, given that the AC method is the least computationally demanding, it seems to be best suited for practical applications. As

Table I: Angular error / coherence / detection rate of the signals obtained using RIRs recorded by Zylia SMA.

HOA	Method	$T_{20} = 1.21$ s	$T_{20} = 0.77$ s	$T_{20} = 0.57$ s	$T_{20} = 0.45$ s	$T_{20} = 0.37$ s
1	AC	<b>44.75° / 0.71 / 47%</b>	40.27° / 0.81 / 47%	39.53° / <b>0.83</b> / 47%	35.32° / 0.83 / 47%	32.35° / 0.82 / 47%
	COV	44.82° / <b>0.71</b> / 47%	<b>38.72° / 0.83</b> / 47%	<b>38.04°</b> / 0.81 / 47%	<b>33.79° / 0.83</b> / 53%	<b>30.80° / 0.85</b> / 47%
	ADMM	46.69° / 0.70 / 47%	40.75° / 0.81 / 47%	44.14° / 0.78 / <b>53%</b>	40.46° / 0.78 / 53%	38.28° / 0.80 / <b>53%</b>
	MCLMS	84.82° / 0.31 / 7%	85.95° / 0.30 / 7%	84.66° / 0.32 / 7%	80.4° / 0.35 / 7%	84.32° / 0.30 / 7%
	RNMCFLMS	85.00° / 0.26 / 27%	83.08° / 0.28 / 27%	81.57° / 0.38 / 27%	82.19° / 0.28 / 20%	88.92° / 0.26 / 20%
	GTVV	59.11° / 0.55 / 40%	61.43° / 0.63 / 40%	65.50° / 0.57 / 47%	61.37° / 0.63 / 47%	60.10° / 0.47 / 47%
	TDVV	73.4° / 0.30 / 33%	67.85° / 0.44 / 33%	74.29° / 0.29 / 40%	70.89° / 0.23 / 40%	63.86° / 0.35 / 40%
2	AC	<b>33.72° / 0.58</b> / 47%	<b>21.30° / 0.77</b> / 47%	<b>25.11° / 0.72</b> / 53%	26.33° / 0.73 / <b>53%</b>	17.63° / 0.76 / <b>53%</b>
	COV	37.48° / 0.57 / <b>53%</b>	21.82° / 0.75 / <b>53%</b>	26.17° / 0.69 / <b>53%</b>	25.55° / 0.69 / <b>53%</b>	<b>17.62°</b> / 0.78 / <b>53%</b>
	ADMM	<b>33.72°</b> / 0.57 / 47%	22.9° / <b>0.77</b> / 53%	28.36° / <b>0.72</b> / 53%	<b>24.68° / 0.76</b> / 53%	18.69° / <b>0.79</b> / 53%
	MCLMS	83.23° / 0.15 / 7%	84.32° / 0.15 / 7%	77.11° / 0.20 / 7%	73.67° / 0.21 / 7%	74.01° / 0.18 / 7%
	RNMCFLMS	101.93° / -0.07 / 27%	101.93° / -0.02 / 27%	100.7° / -0.01 / 33%	98.78° / 0.10 / 33%	102.83° / 0.04 / 33%
	GTVV	56.07° / 0.41 / 47%	40.14° / 0.56 / 40%	47.87° / 0.54 / 47%	44.69° / 0.62 / <b>53%</b>	40.27° / 0.54 / 47%
	TDVV	67.13° / 0.25 / 40%	62.60° / 0.36 / 40%	72.56° / 0.17 / 40%	73.20° / 0.18 / 47%	59.03° / 0.27 / 47%
3	AC	17.32° / 0.65 / <b>53%</b>	<b>11.02° / 0.80</b> / 53%	<b>13.56° / 0.77</b> / 60%	13.46° / 0.77 / <b>60%</b>	11.58° / 0.79 / <b>60%</b>
	COV	16.92° / <b>0.66</b> / 53%	11.29° / 0.77 / <b>53%</b>	14.34° / 0.74 / 53%	14.77° / 0.74 / <b>60%</b>	11.34° / 0.78 / <b>60%</b>
	ADMM	<b>14.25° / 0.66</b> / 53%	11.40° / <b>0.80</b> / 53%	14.50° / <b>0.77</b> / 60%	<b>12.53° / 0.82</b> / 60%	<b>8.16° / 0.83</b> / 60%
	MCLMS	80.22° / 0.13 / 7%	83.16° / 0.10 / 7%	76.9° / 0.13 / 7%	70.84° / 0.12 / 7%	69.32° / 0.13 / 7%
	RNMCFLMS	98.86° / 0.04 / 27%	98.86° / 0.03 / 27%	98.86° / 0.11 / 33%	99.97° / 0.10 / 33%	101.21° / 0.06 / 33%
	GTVV	39.49° / 0.47 / 47%	21.62° / 0.61 / 47%	31.75° / 0.59 / 47%	18.47° / 0.69 / 53%	20.82° / 0.62 / 53%
	TDVV	68.41° / 0.23 / 40%	64.22° / 0.31 / 40%	66.92° / 0.20 / 47%	66.77° / 0.19 / 47%	56.69° / 0.24 / 47%

Table II: Angular error / coherence / detection rate of the signals obtained using RIRs recorded by Eigenmike SMA.

HOA	Method	$T_{20} = 1.21$ s	$T_{20} = 0.77$ s	$T_{20} = 0.57$ s	$T_{20} = 0.45$ s	$T_{20} = 0.37$ s
1	AC	19.71° / <b>0.91</b> / 53%	19.71° / <b>0.90</b> / 47%	<b>19.88° / 0.90</b> / 53%	18.38° / <b>0.92</b> / 53%	19.99° / <b>0.92</b> / 53%
	COV	<b>19.57° / 0.91</b> / 47%	<b>19.29° / 0.90</b> / 47%	20.24° / <b>0.90</b> / 47%	<b>18.35° / 0.92</b> / 47%	<b>19.85°</b> / 0.91 / <b>53%</b>
	ADMM	19.58° / 0.89 / <b>53%</b>	21.06° / 0.87 / <b>47%</b>	20.64° / <b>0.90</b> / 47%	20.63° / 0.89 / 47%	19.91° / 0.90 / <b>53%</b>
	MCLMS	84.97° / 0.30 / 7%	82.66° / 0.32 / 7%	79.99° / 0.36 / 7%	81.26° / 0.31 / 7%	78.87° / 0.36 / 7%
	RNMCFLMS	73.75° / 0.40 / 33%	69.41° / 0.42 / 33%	61.76° / 0.58 / 27%	68.05° / 0.50 / 27%	57.21° / 0.60 / 27%
	GTVV	39.09° / 0.78 / 47%	45.46° / 0.69 / 40%	40.43° / 0.74 / 47%	37.42° / 0.78 / 53%	36.34° / 0.80 / 47%
	TDVV	46.29° / 0.51 / 47%	48.57° / 0.49 / 40%	50.68° / 0.46 / 47%	49.92° / 0.49 / 40%	46.84° / 0.52 / 47%
2	AC	7.39° / 0.90 / <b>60%</b>	<b>7.77° / 0.88</b> / 53%	8.01° / <b>0.89</b> / 60%	8.21° / <b>0.90</b> / 60%	11.39° / <b>0.89</b> / 60%
	COV	7.39° / 0.90 / 53%	7.91° / <b>0.88</b> / 53%	8.01° / 0.88 / 53%	8.21° / <b>0.90</b> / 60%	11.7° / <b>0.89</b> / 53%
	ADMM	<b>7.23° / 0.91</b> / 53%	<b>7.77°</b> / 0.86 / 53%	<b>7.53° / 0.89</b> / 53%	<b>7.67°</b> / 0.89 / 53%	<b>8.16°</b> / 0.88 / 53%
	MCLMS	80.35° / 0.09 / 7%	84.04° / 0.10 / 7%	70.14° / 0.17 / 7%	79.93° / 0.16 / 7%	64.91° / 0.26 / 7%
	RNMCFLMS	80.54° / 0.13 / 40%	79.50° / 0.09 / 40%	75.43° / 0.17 / 40%	73.37° / 0.20 / 40%	71.56° / 0.24 / 40%
	GTVV	13.02° / 0.84 / 53%	15.29° / 0.79 / <b>53%</b>	13.25° / 0.83 / 53%	13.28° / 0.83 / 53%	13.95° / 0.84 / 53%
	TDVV	35.36° / 0.44 / 47%	39.87° / 0.38 / 47%	33.85° / 0.43 / 53%	38.62° / 0.45 / 47%	38.68° / 0.47 / 53%
3	AC	7.05° / 0.89 / 53%	<b>7.11° / 0.87</b> / 53%	<b>7.05° / 0.87</b> / 60%	<b>7.12° / 0.87</b> / 60%	7.31° / 0.85 / <b>60%</b>
	COV	7.11° / 0.89 / 53%	7.19° / 0.86 / <b>53%</b>	7.12° / 0.85 / 53%	7.20° / 0.86 / <b>60%</b>	7.31° / 0.85 / <b>60%</b>
	ADMM	<b>7.02° / 0.90</b> / 53%	7.19° / 0.86 / <b>53%</b>	7.20° / <b>0.87</b> / 53%	7.23° / <b>0.87</b> / 60%	<b>7.20° / 0.87</b> / 60%
	MCLMS	82.10° / 0.07 / 7%	78.31° / 0.09 / 7%	72.41° / 0.12 / 7%	73.13° / 0.11 / 7%	66.27° / 0.17 / 7%
	RNMCFLMS	86.34° / 0.12 / 33%	84.39° / 0.10 / 33%	80.95° / 0.09 / 40%	82.09° / 0.13 / 40%	79.84° / 0.15 / 40%
	GTVV	7.47° / 0.85 / <b>60%</b>	7.77° / 0.81 / <b>53%</b>	7.77° / 0.82 / <b>60%</b>	8.21° / 0.80 / 53%	7.62° / 0.83 / 53%
	TDVV	30.6° / 0.36 / 47%	39.62° / 0.31 / 47%	30.17° / 0.31 / 53%	31.18° / 0.35 / 53%	32.95° / 0.38 / 53%

expected, amongst baseline approaches, the GTVV representation produces the best results, especially for higher HOA orders. One may also remark that the overall performance of all tested methods improves with the HOA order, and – somewhat surprisingly – is not much affected by the change in sound absorption, *i.e.*, by the RT of the room.

It is important to indicate that some degradation in performance (particularly, in detection rate) may be due to the chosen experimentation protocol, based on pre-selected peaks of the ground truth RdRIR. Related to that, note that while the angular precision and coherence are generally correlated, this is not always the case, suggesting that a more refined method for extracting directions from SH vectors (such as [87]) may further reduce angular errors. This may also explain, to some extent, the disparity between the results obtained from the Eigenmike and Zylia SMAs.

## VI. CONCLUSION

We have presented a detailed discussion on GTVV – Generalized Time-domain Velocity Vector – and proposed several methods for the blind identification of early room impulse responses by exploiting properties of this signal representation. We term the time series extracted by these methods *RdRIR - Reduced Room Impulse Response*. The numerical experiments using simulated and recorded RIRs (acquired by different SMAs) demonstrate the performance gains of RdRIR over the baseline BSI approaches. We envision that some of the proposed techniques, due to their implementation simplicity and small computational overhead, could find their place in many practical applications involving Ambisonics and immersive sound. Future work will focus on improving the angular precision of estimated wavefronts, use of RdRIR representation in learned models (*e.g.* deep neural networks), support for multiple sound sources, and potentially, on extending the

benefits of RdRIR beyond Ambisonics, or even the spatial audio context itself.

## APPENDIX A

### RELATION WITH PSEUDOINTENSITY VECTOR

Sound intensity is defined as the product of acoustic pressure and particle velocity [88], [1]. In a pure-sound field [88], its real part – *active* sound intensity – is orthogonal to the incoming wavefront, suggesting it can be used for determining DoA. The linearized fluid momentum equation states that particle velocity is aligned with the spatial gradient of acoustic pressure [52], hence one needs only an estimate of the acoustic pressure and its gradient to approximate this quantity. Pseudointensity vector is an FOA approximation of active sound intensity, defined as [53], [52]

$$\hat{\mathbf{i}}(f) = \Re \left( \hat{b}_0(f) * \hat{\mathbf{b}}_{1:3}(f) \right), \quad (38)$$

where  $\Re$  denotes the real part of a complex number,  $\hat{b}_0(f)$  is the first (omnidirectional) channel, while  $\hat{\mathbf{b}}_{1:3}(f)$  is the vector of the remaining three FOA channels. Indeed, while  $\hat{b}_0(f)$  is a good approximation of acoustic pressure at the center of an array [66], the other FOA channels exhibit spatial response similar to figure-of-eight microphones aligned with Cartesian coordinate axes [33]. Hence,  $\hat{\mathbf{b}}_{1:3}$  is a decent approximation of the spatial gradient vector.

Note that, for the trivial beamformer  $\mathbf{w} = [1 \ 0 \ 0 \ \dots \ 0]^\top$ , we can express  $\hat{\mathbf{i}}(f)$  using GFVV (3), as follows:

$$\hat{\mathbf{i}}(f) = \frac{1}{|\hat{b}_0(f)|^2} \Re \left( \begin{bmatrix} \hat{v}_1(f) \\ \hat{v}_2(f) \\ \hat{v}_3(f) \end{bmatrix} \right). \quad (39)$$

Thus,  $\hat{\mathbf{i}}(f)$  is parallel to the real part of the RTF vector, excluding its first entry (which is trivially equal to 1).

Consider a very simple scenario: in addition to the wavefront coming from the DoA direction  $(\theta_0, \phi_0)$ , there is an impinging wavefront from a reflected sound in the direction  $(\theta_1, \phi_1)$ . Given the unit response of the omnidirectional channel in all directions, from (3) we can rewrite (39) as

$$\hat{\mathbf{i}}(f) \propto \Re \left( \frac{\vec{u}_0 - \vec{u}_1 \gamma_1}{1 - \gamma_1} \right), \quad (40)$$

where  $\gamma_1 = -\hat{g}_1(f)e^{-j2\pi f\tau_1}$ , while we use  $\vec{u}_0$  and  $\vec{u}_1$  to designate the subvectors composed of entries  $[y[1], y[2], y[3]]^\top$  of the SH encoding vectors  $\mathbf{y}_0$  and  $\mathbf{y}_1$ , respectively.

Disregarding the scaling factors, we have

$$\hat{\mathbf{i}}(f) \propto \vec{u}_0 (1 + \hat{g}_1(f) \cos \varphi_1) + \vec{u}_1 \hat{g}_1(f) (\hat{g}_1(f) + \cos \varphi_1), \quad (41)$$

where  $\varphi_1 = 2\pi f\tau_1$ . Since  $\varphi_1$  varies linearly along frequencies  $f$ , we generally have  $\cos(2\pi f\tau_1) \neq -\hat{g}_1(f)$ , hence the pseudointensity vector  $\hat{\mathbf{i}}(f)$  produces a biased estimate of the DoA direction  $\vec{u}_0$ . At the same time, due to (18) and the assumed  $\hat{g}_1(f) < 1$ , the DoA estimate from GTVV (or even TDVV) representation, would remain unbiased.

## APPENDIX B

### ADMM FOR THE PROBLEM (37)

With  $\mu > 0$ ,  $\tilde{\mathbf{H}}^{(0)} = \mathbf{U}^{(0)} = \mathbf{0}$ , and  $\mathbf{0} \in \mathbb{R}^{(L+1)^2 \times J}$  the all-zero matrix, the proposed ADMM iterates the next steps:

$$\mathbf{H}^{(q+1)} = \underset{\mathbf{H} \in \Xi}{\operatorname{argmin}} \mu \|\mathbf{H}\|_{2,1} + \frac{1}{2} \|\mathbf{H} - \tilde{\mathbf{H}}^{(q)} - \mathbf{U}^{(q)}\|_F^2 \quad (42)$$

$$\mathbf{a}^{(q+1)} = \underset{\mathbf{a}, a_0=1}{\operatorname{argmin}} \sum_{l=0}^{(L+1)^2-1} \|\mathbf{v}_{l,:} * \mathbf{a} - \mathbf{h}_{l,:}^{(q+1)} + \mathbf{u}_{l,:}^{(q)}\|_2^2 \quad (43)$$

$$\tilde{\mathbf{h}}_{l,:}^{(q+1)} = \mathbf{v}_{l,:} * \mathbf{a}^{(q+1)} \quad (44)$$

$$\mathbf{U}^{(q+1)} = \mathbf{U}^{(q)} + \tilde{\mathbf{H}}^{(q+1)} - \mathbf{H}^{(q+1)}, \quad (45)$$

where  $\tilde{\mathbf{h}}_{l,:}$  and  $\mathbf{u}_{l,:}$  denote the  $l^{\text{th}}$  rows of the matrices  $\tilde{\mathbf{H}}$  and  $\mathbf{U}$ , respectively, while  $\Xi$  is the set of all real  $(L+1)^2 \times J$  matrices for which first row has non-negative entries, while the columns  $\mathbf{h}_{:,j}$  indexed by  $j \notin [0, j_{\max}]$  contain only zeros. The mixed norm  $\|\cdot\|_{2,1}$  is equal to the sum of the  $\ell_2$ -norms of matrix columns, while  $\|\mathbf{H}\|_F$  denotes the Frobenius norm of a matrix  $\mathbf{H}$ .

Let  $\mathcal{P}_\Xi(\mathbf{H})$  denote the operator that projects a matrix  $\mathbf{H}$  to  $\Xi$ , *i.e.*, sets to zero all  $h_{0,vj} < 0$  and  $\mathbf{h}_{:,j}$ ,  $j \notin [0, j_{\max}]$ . Define a group soft-thresholding [89] operator  $\mathcal{S}_\mu(\cdot)$  as

$$\mathcal{S}_\mu(\mathbf{H})_{l,j} = \max \left( 0, 1 - \frac{\mu}{\|\mathbf{h}_{:,j}\|_2} \right) h_{l,j}, \quad (46)$$

where  $h_{l,j}$  is an entry of the matrix  $\mathbf{H}$  at the row  $l$  and the column  $j$ . Then, the solution of the subproblem (42) is

$$\mathbf{H}^{(q+1)} = \mathcal{S}_\mu \left( \mathcal{P}_\Xi \left( \tilde{\mathbf{H}}^{(q)} + \mathbf{U}^{(q)} \right) \right). \quad (47)$$

The time complexity of the above operations is linear, *i.e.* of the order  $O((L+1)^2 J)$ .

The subproblem (43) is very similar to the initial constrained quadratic problem (29), and can be rewritten as:

$$\min_{\mathbf{a}} \sum_{l=0}^{(L+1)^2-1} \sum_j ((v_{l,:} * \mathbf{a})_j - d_{l,j})^2, \quad \text{s.t. } a_0 = 1, \quad (48)$$

where  $d_{l,j} = h_{l,j}^{(q+1)} - u_{l,j}^{(q)}$ . It can again be cast into a linear system, similar to (30):

$$\sum_{j=1}^{j_{\max}} a_j r(j, s) = r_{dv}(0, s) - r(0, s), \quad (49)$$

where  $r(j, s)$  are defined as in (31), while the coefficients  $r_{dv}(0, s)$  correspond to the cross-correlation

$$r_{dv}(0, s) = \sum_{l=0}^{(L+1)^2-1} \sum_{j'} d_{l,j'} v_{l,j'-s}. \quad (50)$$

In both cases, the summation with respect to  $j'$  is now done over the entire range  $[-J/2 + 1, -J/2]$ . This means that the autocorrelation coefficients are directly obtained from the power spectrum of GFVV, *i.e.*,

$$r(j, s) = \mathcal{F}^{-1} \left( \sum_{l=0}^{(L+1)^2-1} |\hat{\mathbf{v}}_{l,:}|^2 \right)_{j-s}. \quad (51)$$

Since  $r(j, s)$  is constant across iterations, the Toeplitz matrix and the autocorrelation part of the right hand side of (49) need to be calculated only once.

Analogous to autocorrelation, the cross-correlation values (50) can be computed in frequency domain from the cross-spectra of the involved quantities at  $O((L+1)^2 J \log J)$  cost:

$$r_{dv}(0, s) = \mathcal{F}^{-1} \left( \sum_{l=0}^{(L+1)^2-1} \hat{d}_{l,:} \hat{v}_{l,:}^* \right)_s. \quad (52)$$

Even though we have  $j_{\max} < J/2$  (cf. the discussion in subsection III-B),  $j_{\max}$  and  $J$  are still comparable, hence the per-iteration complexity of the ADMM algorithm is determined by the cost of solving the linear system (49), which requires  $O((j_{\max} + 1)^2)$  operations.

In practice, the convergence speed depends on the parameter  $\mu$ , which is set to 0.1 in our experiments. Convergence criterion based on the primal and dual updates has been presented in [76], however, we observe that the algorithm typically produces meaningful results within tens of iterations. Furthermore, the algorithm can be accelerated by warm-starting, *i.e.*, by initializing the iterations with the estimates from the previous frame.

## REFERENCES

- [1] H. Kuttruff, *Room acoustics*, CRC Press, 2016.
- [2] C. Evers, H. Löllmann, H. Mellmann, A. Schmidt, H. Barfuss, P. Naylor, and W. Kellermann, "The LOCATA challenge: Acoustic source localization and tracking," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1620–1643, 2020.
- [3] K. Kinoshita, M. Delcroix, S. Gannot, E. Habets, R. Haeb-Umbach, W. Kellermann, V. Leutnant, R. Maas, T. Nakatani, B. Raj, et al., "A summary of the reverb challenge: state-of-the-art and remaining challenges in reverberant speech processing research," *EURASIP Journal on Advances in Signal Processing*, vol. 2016, no. 1, pp. 1–19, 2016.
- [4] J. Zhang, C. Zorilă, R. Doddipatla, and J. Barker, "On end-to-end multi-channel time domain speech separation in reverberant environments," in *ICASSP 2020 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2020, pp. 6389–6393.
- [5] A. O'Donovan, R. Duraiswami, and D. Zotkin, "Automatic matched filter recovery via the audio camera," in *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2010, pp. 2826–2829.
- [6] D. Di Carlo, A. Deleforge, and N. Bertin, "Mirage: 2d source localization using microphone pair augmentation with echoes," in *ICASSP 2019 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2019, pp. 775–779.
- [7] J. Daniel and S. Kitić, "Time domain velocity vector for retracing the multipath propagation," in *ICASSP 2020 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2020, pp. 421–425.
- [8] R. Scheibler, D. Di Carlo, A. Deleforge, and I. Dokmanic, "Separake: Source separation with a little help from echoes," in *ICASSP 2018 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2018, pp. 6897–6901.
- [9] R. Weisman, T. Shlomo, V. Tourbabin, P. Calamia, and B. Rafaely, "Robustness of Acoustic Rake Filters in Minimum Variance Beamforming," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3668–3678, 2021.
- [10] O. Shmaryahu and S. Gannot, "On the importance of acoustic reflections in beamforming," in *IWAENC 2022 - International Workshop on Acoustic Signal Enhancement*. IEEE, 2022, pp. 1–5.
- [11] R. Giri, M. Seltzer, J. Droppo, and D. Yu, "Improving speech recognition in reverberation using a room-aware deep neural network and multi-task learning," in *ICASSP 2015 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2015, pp. 5014–5018.
- [12] M. Yasuda, Y. Ohishi, and S. Saito, "Echo-aware adaptation of sound event localization and detection in unknown environments," in *ICASSP 2022 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2022, pp. 226–230.
- [13] A. Ratnarajah, I. Ananthabhotla, V. K. Ithapu, P. Hoffmann, D. Manocha, and P. Calamia, "Towards improved room impulse response estimation for speech recognition," in *ICASSP 2023 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2023, pp. 1–5.
- [14] S. Kitić, N. Bertin, and R. Gribonval, "Hearing behind walls: localizing sources in the room next door with cosparsity," in *ICASSP 2014 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2014, pp. 3087–3091.
- [15] I. An, M. Son, D. Manocha, and S.-E. Yoon, "Reflection-aware sound source localization," in *ICRA 2018 - IEEE International Conference on Robotics and Automation*. IEEE, 2018, pp. 66–73.
- [16] J. Boger-Lombard, Y. Slobodkin, and O. Katz, "Towards passive non-line-of-sight acoustic localization around corners using uncontrolled random noise sources," *Scientific Reports*, vol. 13, no. 1, pp. 4952, 2023.
- [17] I. Dokmanić, R. Parhizkar, A. Walther, Y. Lu, and M. Vetterli, "Acoustic echoes reveal room shape," *Proceedings of the National Academy of Sciences*, vol. 110, no. 30, pp. 12186–12191, 2013.
- [18] M. Lovedee-Turner and D. Murphy, "Three-dimensional reflector localisation and room geometry estimation using a spherical microphone array," *The Journal of the Acoustical Society of America*, vol. 146, no. 5, pp. 3339–3352, 2019.
- [19] F. Antonacci, J. Filos, M. Thomas, E. Habets, A. Sarti, P. Naylor, and S. Tubaro, "Inference of room geometry from acoustic impulse responses," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 10, pp. 2683–2695, 2012.
- [20] F. Ribeiro, C. Zhang, D. Florêncio, and D. Ba, "Using reverberation to improve range and elevation discrimination for small array sound source localization," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 7, pp. 1781–1792, 2010.
- [21] L. Birnie, T. Abhayapala, and P. Samarasinghe, "Reflection assisted sound source localization through a harmonic domain music framework," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 279–293, 2019.
- [22] J. Daniel and S. Kitić, "Echo-enabled Direction-of-Arrival and range estimation of a mobile source in Ambisonic domain," in *EUSPICO 2022 - 30th European Signal Processing Conference*. IEEE, 2022, pp. 852–856.
- [23] W. Yu and B. Kleijn, "Room acoustical parameter estimation from room impulse responses using deep neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 436–447, 2020.
- [24] S. Dilungana, A. Deleforge, C. Foy, and S. Faisan, "Geometry-informed estimation of surface absorption profiles from room impulse responses," in *EUSPICO 2022 - 30th European Signal Processing Conference*. EURASIP, 2022.
- [25] M. Baum, L. Cuccovillo, A. Yaroshchuk, and P. Aichroth, "Environment classification via blind roomprints estimation," in *WIFS 2022 - IEEE International Workshop on Information Forensics and Security*. IEEE, 2022, pp. 1–6.
- [26] J. Su, Z. Jin, and A. Finkelstein, "Acoustic matching by embedding impulse responses," in *ICASSP 2020 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2020, pp. 426–430.
- [27] Z. Tang, N. Bryan, D. Li, T. Langlois, and D. Manocha, "Scene-aware audio rendering via deep acoustic analysis," *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, no. 5, pp. 1991–2001, 2020.
- [28] C.-Y. Chi, C.-C. Feng, C.-H. Chen, and C.-Y. Chen, *Blind equalization and system identification: batch processing algorithms, performance and applications*, Springer Science & Business Media, 2006.
- [29] G. Xu, H. Liu, L. Tong, and T. Kailath, "A least-squares approach to blind channel identification," *IEEE Transactions on Signal Processing*, vol. 43, no. 12, pp. 2982–2993, 1995.
- [30] P. Naylor, N. Gaubitch, et al., *Speech dereverberation*, vol. 2, Springer, 2010.
- [31] S. Kitić and J. Daniel, "Generalized Time Domain Velocity Vector," in *ICASSP 2022 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2022, pp. 936–940.
- [32] D. Jarrett, E. Habets, and P. Naylor, *Theory and applications of spherical microphone array processing*, vol. 9, Springer, 2017.

- [33] F. Zotter and M. Frank, *Ambisonics: A practical 3D audio theory for recording, studio production, sound reinforcement, and virtual reality*, Springer Nature, 2019.
- [34] H. Lee, "Multichannel 3d microphone arrays: A review," *Journal of the Audio Engineering Society*, vol. 69, no. 1/2, pp. 5–26, 2021.
- [35] Y. Sato, "A method of self-recovering equalization for multilevel amplitude-modulation systems," *IEEE Transactions on communications*, vol. 23, no. 6, pp. 679–682, 1975.
- [36] L. Tong, G. Xu, and T. Kailath, "A new approach to blind identification and equalization of multipath channels," in *Conference Record of the Twenty-Fifth Asilomar Conference on Signals, Systems & Computers*. IEEE Computer Society, 1991, pp. 856–857.
- [37] E. Moulines, P. Duhamel, J.-F. Cardoso, and S. Mayrargue, "Subspace methods for the blind identification of multichannel fir filters," *IEEE Transactions on Signal Processing*, vol. 43, no. 2, pp. 516–525, 1995.
- [38] Y. Hua, "Fast maximum likelihood for blind identification of multiple fir channels," *IEEE Transactions on Signal Processing*, vol. 44, no. 3, pp. 661–672, 1996.
- [39] J. Benesty, "Adaptive eigenvalue decomposition algorithm for passive acoustic source localization," *The Journal of the Acoustical Society of America*, vol. 107, no. 1, pp. 384–391, 2000.
- [40] Y. Huang and J. Benesty, "Adaptive multi-channel least mean square and newton algorithms for blind channel identification," *Signal processing*, vol. 82, no. 8, pp. 1127–1138, 2002.
- [41] Y. Huang and J. Benesty, "A class of frequency-domain adaptive approaches to blind multichannel identification," *IEEE Transactions on Signal Processing*, vol. 51, no. 1, pp. 11–24, 2003.
- [42] R. Ahmad, A. Khong, and P. Naylor, "Proportionate frequency domain adaptive algorithms for blind channel identification," in *ICASSP 2006 - IEEE International Conference on Acoustics Speech and Signal processing Proceedings*. IEEE, 2006, vol. 5, pp. V–V.
- [43] W. Xue, M. Brookes, and P. Naylor, "Cross-correlation based under-modelled multichannel blind acoustic system identification with sparsity regularization," in *EUSIPCO 2016 - 24th European Signal Processing Conference*. IEEE, 2016, pp. 718–722.
- [44] W. Xue, M. Brookes, and P. Naylor, "Frequency-domain under-modelled blind system identification based on cross power spectrum and sparsity regularization," in *ICASSP 2017 - IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2017, pp. 591–595.
- [45] S. Wager, K. Choi, and S. Durand, "Dereverberation using joint estimation of dry speech signal and acoustic system," *arXiv preprint arXiv:2007.12581*, 2020.
- [46] C. Steinmetz, V. K. Ithapu, and P. Calamia, "Filtered noise shaping for time domain room impulse response estimation from reverberant speech," in *WASPAA 2021 - IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*. IEEE, 2021, pp. 221–225.
- [47] M. Pezzoli, D. Perini, A. Bernardini, F. Borra, F. Antonacci, and A. Sarti, "Deep prior approach for room impulse response reconstruction," *Sensors*, vol. 22, no. 7, pp. 2710, 2020.
- [48] S. Gannot, E. Vincent, S. Markovich-Golan, and A. Ozerov, "A consolidated perspective on multimicrophone speech enhancement and source separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 4, pp. 692–730, 2017.
- [49] L. Perotin, R. Serizel, E. Vincent, and A. Guérin, "Multichannel speech separation with recurrent neural networks from high-order ambisonics recordings," in *ICASSP 2018 - IEEE International Conference on Acoustics, Speech and Signal processing*. IEEE, 2018, pp. 36–40.
- [50] A. Bosca, A. Guérin, L. Perotin, and S. Kitić, "Dilated U-net based approach for multichannel speech enhancement from First-Order Ambisonics recordings," in *EUSIPCO 2020 - 28th European Signal Processing Conference*. IEEE, 2021, pp. 216–220.
- [51] Y. Hu, P. Samarasinghe, S. Gannot, and T. Abhayapala, "Semi-supervised multiple source localization using relative harmonic coefficients under noisy and reverberant environments," *IEEE/ACM Transactions on Audio, Speech, and Language processing*, vol. 28, pp. 3108–3123, 2020.
- [52] J. Merimaa, *Analysis, synthesis, and perception of spatial sound: binaural localization modeling and multichannel loudspeaker reproduction*, Ph.D. thesis, Helsinki University of Technology, 2006.
- [53] D. Jarrett, E. Habets, and P. Naylor, "3d source localization in the spherical harmonic domain using a pseudointensity vector," in *EUSIPCO 2010 - 18th European Signal processing Conference*. IEEE, 2010, pp. 442–446.
- [54] S. Meier and W. Kellermann, "Analysis of the performance and limitations of ica-based relative impulse response identification," in *EUSIPCO 2015 - 23rd European Signal Processing Conference*. IEEE, 2015, pp. 414–418.
- [55] S. Gannot, D. Burshtein, and E. Weinstein, "Signal enhancement using beamforming and nonstationarity with applications to speech," *IEEE Transactions on Signal Processing*, vol. 49, no. 8, pp. 1614–1626, 2001.
- [56] L. Gölles and F. Zotter, "Directional enhancement of first-order ambisonic room impulse responses by the 2+2 directional signal estimator," in *Proceedings of the 15th International Conference on Audio Mostly*, 2020, pp. 38–45.
- [57] A. Herzog and E. Habets, "Generalized intensity vector and energy density in the spherical harmonic domain: Theory and applications," *The Journal of the Acoustical Society of America*, vol. 150, no. 1, pp. 294–306, 2021.
- [58] C. Borrelli, A. Canclini, F. Antonacci, A. Sarti, and S. Tubaro, "A denoising methodology for higher order ambisonics recordings," in *IWAENC 2018 - 16th International Workshop on Acoustic Signal Enhancement*. IEEE, 2018, pp. 451–455.
- [59] N. Meyer-Kahlen and S. Schlecht, "Blind directional room impulse response parameterization from relative transfer functions," in *IWAENC 2022 - International Workshop on Acoustic Signal Enhancement*. IEEE, 2022, pp. 1–5.
- [60] L. Madmoni and B. Rafaely, "Direction of arrival estimation for reverberant speech based on enhanced decomposition of the direct sound," *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 1, pp. 131–142, 2018.
- [61] Y. Biderman, B. Rafaely, S. Gannot, and S. Doclo, "Efficient relative transfer function estimation framework in the spherical harmonics domain," in *EUSIPCO 2016 - 24th European Signal Processing Conference*. IEEE, 2016, pp. 1658–1662.
- [62] T. Shlomo and B. Rafaely, "Blind localization of early room reflections using phase aligned spatial correlation," *IEEE Transactions on Signal Processing*, vol. 69, pp. 1213–1225, 2021.
- [63] T. Hidaka, Y. Yamada, and T. Nakagawa, "A new definition of boundary point between early reflections and late reverberation in room impulse responses," *The Journal of the Acoustical Society of America*, vol. 122, no. 1, pp. 326–332, 2007.
- [64] J. Allen and D. Berkley, "Image method for efficiently simulating small-room acoustics," *The Journal of the Acoustical Society of America*, vol. 65, no. 4, pp. 943–950, 1979.
- [65] W. Rudin, "Functional analysis 2nd ed," *International Series in Pure and Applied Mathematics*. McGraw-Hill, Inc., New York, 1991.
- [66] B. Rafaely, *Fundamentals of spherical array processing*, vol. 8, Springer, 2015.
- [67] O. Shalvi and E. Weinstein, "System identification using nonstationary signals," *IEEE Transactions on Signal Processing*, vol. 44, no. 8, pp. 2055–2063, 1996.
- [68] S. Markovich-Golan, S. Gannot, and W. Kellermann, "Performance analysis of the covariance-whitening and the covariance-subtraction methods for estimating the relative transfer function," in *EUSIPCO 2018 - 26th European Signal Processing Conference*. IEEE, 2018, pp. 2499–2503.
- [69] D. Cassioli and A. Meocozzi, "Minimum-phase impulse response channels," *IEEE Transactions on communications*, vol. 57, no. 12, pp. 3529–3532, 2009.
- [70] G. Golub and C. Van Loan, *Matrix computations*, JHU press, 2013.
- [71] J. Daniel and S. Moreau, "Further study of sound field coding with Higher Order Ambisonics," in *Audio Engineering Society Convention 116*. Audio Engineering Society, 2004.
- [72] M. Hayes, *Statistical digital signal processing and modeling*, John Wiley & Sons, 2009.
- [73] K. Steiglitz and L. McBride, "A technique for the identification of linear systems," *IEEE Transactions on Automatic Control*, vol. 10, no. 4, pp. 461–464, 1965.
- [74] J. Eckstein and D. Bertsekas, "On the Douglas—Rachford splitting method and the proximal point algorithm for maximal monotone operators," *Mathematical Programming*, vol. 55, no. 1, pp. 293–318, 1992.
- [75] A. Aragón, J. Francisco, R. Campoy, and M. Tam, "The Douglas—Rachford algorithm for convex and nonconvex feasibility problems," *Mathematical Methods of Operations Research*, vol. 91, pp. 201–240, 2020.
- [76] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [77] M. Haque and M. Hasan, "Noise robust multichannel frequency-domain LMS algorithms for blind channel identification," *IEEE Signal Processing Letters*, vol. 15, pp. 305–308, 2008.
- [78] E. Habets and P. Naylor, "Blind System Identification and Equalization Toolbox," 2009.

- [79] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *2015 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2015, pp. 5206–5210.
- [80] J. Daniel, “Spatial Sound Encoding Including Near Field Effect: Introducing Distance Coding Filters and a Viable, New Ambisonic Format,” in *Audio Engineering Society 23rd International Conference*. Audio Engineering Society, 2003.
- [81] T. Gerkmann and R. Hendriks, “Unbiased MMSE-based noise power estimation with low complexity and low tracking delay,” *IEEE Transactions on Audio, Speech, and Lang. Processing*, vol. 20, no. 4, pp. 1383–1393, 2011.
- [82] E. Habets, “Room impulse response generator,” *Technische Universiteit Eindhoven, Tech. Rep.*, vol. 2, no. 2.4, pp. 1, 2006.
- [83] E. Vincent and D. Campbell, “Roomsimove: Matlab toolbox for the computation of simulated room impulse responses for moving sources,” 2015.
- [84] B. Efron, *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*, vol. 1, Cambridge University Press, 2012.
- [85] T. Yoshioka and T. Nakatani, “Generalization of multi-channel linear prediction methods for blind MIMO impulse response shortening,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 10, pp. 2707–2720, 2012.
- [86] T. McKenzie, L. McCormack, and C. Hold, “Dataset of spatial room impulse responses in a variable acoustics room for six degrees-of-freedom rendering and analysis,” *arXiv preprint arXiv:2111.11882*, 2021.
- [87] A. Herzog and E. Habets, “Eigenbeam-ESPRIT for DOA-vector estimation,” *IEEE Signal Processing Letters*, vol. 26, no. 4, pp. 572–576, 2019.
- [88] F. Jacobsen, “A note on instantaneous and time-averaged active and reactive sound intensity,” *Journal of Sound and Vibration*, vol. 147, no. 3, pp. 489–496, 1991.
- [89] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski, “Optimization with sparsity-inducing penalties,” *Foundations and Trends® in Machine Learning*, vol. 4, no. 1, pp. 1–106, 2012.