

Blind Localization Of Early Room Reflections Using Phase Aligned Spatial Correlation

Tom Shlomo, *Student Member, IEEE*, Boaz Rafaely, *Senior Member, IEEE*

Abstract—Blind estimation of the direction of arrival (DOA) and delay of room reflections from reverberant sound may be useful for a wide range of applications. However, due to the high temporal and spatial density of early room reflections and their low power compared to the direct sound, existing methods can only detect a small number of reflections. This paper presents PHALCOR (PHase ALigned CORrelation), a novel method for blind estimation of the DOA and delay of early reflections that overcomes the limitations of existing solutions. PHALCOR is based on a signal model in which the reflection signals are explicitly modeled as delayed and scaled copies of the direct sound. A phase alignment transform of the spatial correlation matrices is proposed; this transform can separate reflections with different delays, enabling the detection and localization of reflections with similar DOAs. It is shown that the DOAs and delays of the early reflections can be estimated by separately analysing the left and right singular vectors of the transformed matrices using sparse recovery techniques. An extensive simulation study of a speaker in a reverberant room, recorded by a spherical array, demonstrates the effectiveness of the proposed method.

Index Terms—Direction-of-arrival estimation, room reflections, MUSIC, Spherical array, sparse recovery

I. INTRODUCTION

Estimation of the direction of arrival (DOA) and delay of room reflections is useful for many tasks in signal processing, such as speech enhancement and dereverberation [1], [2], source separation [3], optimal beamforming [4] and room geometry inference [5]. The early reflections have a key role in sound perception, as they can improve speech intelligibility and listener envelopment. They are also related to the impression of source width, loudness and distance [6], [7]. Therefore, exploitation of the early reflections can be beneficial in parametric spatial audio methods and spatial audio coding [8], [9].

Existing methods for the estimation of the parameters of early reflections can be categorized as blind and non-blind. Non-blind methods, operate on room impulse response signals, or, alternatively, assume that an anechoic recording of the sound source is available. Blind methods operate on microphone signals directly, and assume that no other information is available, as is often the case in many real world applications. This work focuses on blind estimation.

Spatial filtering, i.e. beamforming, can be utilized to blindly estimate the DOAs of the early reflections, as well as to separate reflection signals from the direct sound, which enables

delay estimation using cross-correlation analysis [5]. However, when arrays of practical size are used, the spatial resolution achieved by beamformers is often insufficient, as the spatial density of early reflections can be very high [10]. Subspace methods, such as MUSIC or ESPRIT [11], [12], can often provide higher resolution than beamformers. However, these methods require the sources to be uncorrelated, while in the case of early reflections, since all sources are delayed copies of the direct sound, reflected narrowband signals are highly correlated. Frequency smoothing is a common method to decorrelate source signals, enabling the application of subspace methods. However, frequency smoothing cannot decorrelate reflections that have similar delays. Furthermore, subspace methods typically require an estimation of the number of sources, which is a challenging task when the amplitudes of the sources greatly vary, as in the case of early reflections. Also, these methods require that the number of microphones is larger than the number of significant reflections, which can limit the number of detected reflections. By formulating the problem as an under-determined linear system, sparse recovery can also be utilized for the localization of early reflections [13]. Since sparsity based methods attempt to find the smallest set of sources that explain the measured signals, their performance improves as the actual number of sources is reduced, and in practice only the first few reflections are recoverable with practical arrays. A common limitation of the methods mentioned above is the use of a multiple source model that does not distinguish between sources with the same DOA. In summary, due to the challenging nature of this task, and the limitations of current methods, no adequate solution seems to be available for blind estimation of the DOA and delay of early room reflections.

This paper presents PHALCOR (PHase ALigned CORrelation), a novel method for blind estimation of the DOA and delay of early reflections. The proposed approach utilizes the inherent structure of early reflections - they are delayed and attenuated copies of the direct sound. More specifically, we use the property that the narrowband correlation between a source and its reflections has a phase that is linear in frequency to construct a transform that can separate reflections with different delays, which also enables the detection of multiple reflections from the same direction. Since the number of reflections with similar delays is usually small, the DOAs of reflections with similar delays are estimated using orthogonal matching pursuit (OMP), a sparse recovery technique. A simulation study demonstrates the performance of PHALCOR, in particular its ability to accurately detect a large number of reflections. Initial results of this work, with a simplified method and a much-

The authors are with the School of Electrical and Computer Engineering, Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel (e-mail: tomshlomo@gmail.com; br@bgu.ac.il).

This research was supported by Facebook Reality Labs.

reduced theoretical and performance analysis, were submitted for publication at the Interspeech 2020 conference [14].

The rest of this paper is organized as follows. Section II presents the necessary mathematical background on the plane wave amplitude density function and its spherical Fourier transform. Section III presents the system model. Section IV presents the theoretical foundations of the proposed method, while section V describes the proposed algorithm. A simulation study and conclusions are presented in sections VI and VII, respectively.

II. MATHEMATICAL BACKGROUND

A. Notation

The notation used in the paper is presented briefly in this section. Lowercase boldface letters denote vectors, and uppercase boldface letter denote matrices. The k, l entry of a matrix \mathbf{A} is denoted by $[\mathbf{A}]_{k,l}$. The complex conjugate, transpose, and conjugate transpose are denoted by $(\cdot)^*$, $(\cdot)^T$ and $(\cdot)^H$ respectively. The Euclidean norm of a vector is denoted by $\|\cdot\|$. The outer-product of two vectors \mathbf{a} and \mathbf{b} is the matrix $\mathbf{a}\mathbf{b}^H$. The imaginary unit is denoted by i .

\mathbb{S}^2 denotes the unit sphere in \mathbb{R}^3 . The symbol $\Omega \in \mathbb{S}^2$ represents a direction in 3D-space, i.e a pair of azimuth-elevation angles. $\angle(\Omega, \Omega') \triangleq \arccos(\Omega \cdot \Omega')$ is the angle between directions Ω and Ω' . “ \cdot ” is the dot product in \mathbb{R}^3 .

B. Sound Field Representation Using Plane Wave Amplitude Density

Consider a sound field composed of M plane waves with amplitudes $a_1(f), \dots, a_M(f)$ at frequency f , and directions $\Omega_1, \dots, \Omega_M$. The sound pressure p at any point in space $\mathbf{x} \in \mathbb{R}^3$ can be formulated as follows:

$$p(f, \mathbf{x}) = \sum_{m=1}^M a_m(f) e^{ik\Omega_m \cdot \mathbf{x}} \quad (1)$$

where $k = 2\pi f/c$ is the wave-number, and c is the speed of sound. When the sound field is composed of a continuum of plane waves, the summation is replaced by an integral over the entire sphere, and the amplitudes are replaced by the plane wave amplitude density (PWAD) $a(f, \Omega)$:

$$p(f, \mathbf{x}) = \int_{\mathbb{S}^2} a(f, \Omega) e^{ik\Omega \cdot \mathbf{x}} d\Omega$$

For a fixed frequency, the PWAD is a function on the unit sphere. As such, it is possible to describe it by its spherical Fourier transform (SFT) coefficients [15]:

$$a_{n,m}(f) \triangleq \int_{\mathbb{S}^2} a(f, \Omega) [Y_n^m(\Omega)]^* d\Omega \quad (2)$$

where Y_n^m is the order- n and degree- m spherical harmonic. The SFT of the PWAD can be used to represent the sound pressure as follows [15]:

$$p(f, \mathbf{x}) = \sum_{n=0}^{\infty} \sum_{m=-n}^n 4\pi i^n j_n(kr) Y_n^m(\Omega) a_{n,m}(f) \quad (3)$$

where $r = \|\mathbf{x}\|$, $\Omega = \mathbf{x}/r$, and j_n is the n 'th order spherical Bessel function of the first kind.

Equation (3) can be well approximated by truncating the infinite sum to order $N = \lceil kr \rceil$ [15]. A microphone array can be used to estimate the coefficients of the SFT of the PWAD with order less than or equal to N , by inverting the (truncated) linear transformation (3), a process known as plane wave decomposition. The existence and stability of the inverse transform depend on the frequency f and physical properties of the (typically spherical) microphone array. Further details can be found in [15]. The resulting signals are stacked in a vector of length $(N+1)^2$ as follows:

$$\mathbf{a}_{\text{nm}} \triangleq [a_{0,0}, a_{1,-1}, a_{1,0}, a_{1,1}, \dots, a_{N,N}]^T$$

The rest of this paper is described in terms of the SFT of the PWAD. Processing and analysis in this domain offer several advantages. First, the PWAD provides a description of the sound-field that is independent of the microphone array. Second, the steering vectors, i.e the response to a plane-wave from a given direction, are frequency independent. The steering vector $\mathbf{y}(\Omega)$ is given by [15]:

$$\mathbf{y}(\Omega) \triangleq \frac{\sqrt{4\pi}}{N+1} [Y_0^0(\Omega), Y_1^{-1}(\Omega), \dots, Y_N^N(\Omega)]^H \quad (4)$$

The constant $\frac{\sqrt{4\pi}}{N+1}$ was chosen for convenience such that $\|\mathbf{y}(\Omega)\| = 1$.

III. SYSTEM MODEL

This section presents the system model used in the paper. Consider a sound field comprised of a single source in a room, with a frequency domain signal $s(f)$, and a DOA Ω_0 , relative to a measurement point in the room. As the sound from the source propagates in the room, it is reflected from the room boundaries. The k 'th reflection is modeled as a separate source with DOA Ω_k and signal $s_k(f)$, which is a delayed and scaled copy of the source signal [16]:

$$s_k(f) = \alpha_k e^{-i2\pi f \tau_k} s(f) \quad (5)$$

where τ_k is the delay relative to the direct sound, and α_k is the scaling factor. τ_0 and α_0 are accordingly normalized to 0 and 1, respectively. It is assumed that the delays are sorted such that $\tau_{k-1} \leq \tau_k$.

Let $\mathbf{a}_{\text{nm}}(f)$ denote the vector of the SFT coefficients of the PWAD, up to order N , as a function of frequency. Assuming the sources are in the far field, $\mathbf{a}_{\text{nm}}(f)$ is described by the following model:

$$\mathbf{a}_{\text{nm}}(f) = \mathbf{Y}\mathbf{s}(f) + \mathbf{n}(f) \quad (6)$$

where:

$$\mathbf{s}(f) \triangleq [s_0(f), \dots, s_K(f)]^T \quad (7)$$

$$\mathbf{Y} \triangleq [\mathbf{y}(\Omega_0), \dots, \mathbf{y}(\Omega_K)] \quad (8)$$

$\mathbf{n}(f)$ includes noise and late reverberation terms, and K is the number of early reflections.

Let $\mathbf{R}(f)$ denote the spatial correlation matrix (SCM) at frequency f :

$$\mathbf{R}(f) \triangleq \mathbb{E} [\mathbf{a}_{\text{nm}}(f) \mathbf{a}_{\text{nm}}(f)^H] \quad (9)$$

Substituting Eq. (6) into Eq. (9), and assuming $\mathbf{n}(f)$ and $s(f)$ are uncorrelated, yields:

$$\mathbf{R}(f) = \mathbf{Y}\mathbf{M}(f)\mathbf{Y}^H + \mathbf{N}(f) \quad (10)$$

where:

$$\mathbf{N}(f) \triangleq \mathbb{E} [\mathbf{n}(f)\mathbf{n}(f)^H] \quad (11)$$

$$\mathbf{M}(f) \triangleq \mathbb{E} [\mathbf{s}(f)\mathbf{s}(f)^H] \quad (12)$$

IV. PHASE ALIGNMENT OF THE SPATIAL CORRELATION MATRICES

PHALCOR is based on a phase alignment transformation of the SCM. This section presents the definition and properties of this transformation.

A. Motivation

Before presenting the mathematical details, we first provide some motivation for this transformation.

Equation (10) can be rewritten as:

$$\mathbf{R}(f) = \sum_{k=0}^K \sum_{k'=0}^K [\mathbf{M}(f)]_{k,k'} \mathbf{y}(\Omega_k)\mathbf{y}(\Omega_{k'})^H + \mathbf{N}(f) \quad (13)$$

It is apparent from equation (13) that neglecting \mathbf{N} , the matrix \mathbf{R} is a mixture of the outer products of the steering vectors of the sources. The mixing coefficients are the entries of \mathbf{M} , and therefore it is henceforth referred to as the mixing matrix. Note that the mixing coefficients are frequency dependent, but the steering vectors are not.

Suppose that the k, k' entry of $\mathbf{M}(f)$ has a dominant magnitude, relative to all other entries. This leads to:

$$\mathbf{R}(f) \approx c\mathbf{y}(\Omega_k)\mathbf{y}(\Omega_{k'})^H + \mathbf{N}(f) \quad (14)$$

for some $c \in \mathbb{C}$. Intuitively, estimating Ω_k and $\Omega_{k'}$ in such a case is easier than in the general case. However, according to (12), there may not be a dominant entry in \mathbf{M} . Not only it is Hermitian, but also the magnitudes of its entries are products of amplitudes between pairs of two sources. Assuming the amplitude of the direct sound may be dominant, the $0, 0$ 'th entry that corresponds to the direct sound only may indeed be dominant. However, this is not helpful for localizing the early reflections. The processing presented below is designed to enhance specific entries in \mathbf{M} , so that specific reflections can be more easily localized.

B. Phase aligned Spatial Correlation

We define the following matrix, which we call the phase aligned SCM:

$$\bar{\mathbf{R}}(\tau, f) \triangleq \sum_{j=0}^{J_f-1} w_j \mathbf{R}(f + j\Delta f) e^{i2\pi\tau j\Delta f} \quad (15)$$

where $\tau \geq 0$, Δf is the frequency resolution, and J_f is an integer parameter representing the overall number of frequency points. w_0, \dots, w_{J_f-1} are non-negative weights. Note that when $\tau = 0$ and $w_j = 1$, $\bar{\mathbf{R}}$ is identical to the SCM obtained by frequency smoothing. The matrices $\bar{\mathbf{N}}(\tau, f)$ and $\bar{\mathbf{M}}(\tau, f)$

are similarly defined by replacing \mathbf{R} in (15) with \mathbf{N} and \mathbf{M} , respectively, such that:

$$\bar{\mathbf{R}}(\tau, f) = \sum_{k=0}^K \sum_{k'=0}^K [\bar{\mathbf{M}}(\tau, f)]_{k,k'} \mathbf{y}(\Omega_k)\mathbf{y}(\Omega_{k'})^H + \bar{\mathbf{N}}(\tau, f) \quad (16)$$

Similarly to Eq. (13), Eq. (16) presents $\bar{\mathbf{R}}$ as a mixture of outer-products of pairs of sources' steering vectors. Next, it is proven that for a fixed f , the k, k' entry of $\bar{\mathbf{M}}(\tau, f)$ is maximized for $\tau = \tau_k - \tau_{k'}$. We begin the proof by deriving an explicit expression for the absolute value of the entries of $\bar{\mathbf{M}}$ for an arbitrary τ :

$$\begin{aligned} & \left| [\bar{\mathbf{M}}(\tau, f)]_{k,k'} \right| \\ &= \left| \sum_{j=0}^{J_f-1} [\mathbf{M}(f_j)]_{k,k'} w_j e^{i2\pi\tau j\Delta f} \right| \\ &= \left| \sum_{j=0}^{J_f-1} \mathbb{E} [s_k(f_j)s_{k'}(f_j)] w_j e^{i2\pi\tau j\Delta f} \right| \\ &= \left| \sum_{j=0}^{J_f-1} \alpha_k \alpha_{k'}^* \sigma_s^2(f_j) w_j e^{i2\pi\tau j\Delta f (\tau - (\tau_k - \tau_{k'}))} \right| \end{aligned} \quad (17)$$

where $f_j = f + j\Delta f$ and

$$\sigma_s^2(f_j) \triangleq \mathbb{E} [|s(f_j)|^2] \quad (18)$$

The second equality in Eq. 17 is due to the definition of $\mathbf{M}(f)$ in Eq. (12), while the third equality is due to Eq. (5). Now, by a simple use of the triangle-inequality:

$$\begin{aligned} & \left| [\bar{\mathbf{M}}(\tau, f)]_{k,k'} \right| \leq \sum_{j=0}^{J_f-1} |\alpha_k \alpha_{k'}^* \sigma_s^2(f_j) w_j| \\ &= \left| [\bar{\mathbf{M}}(\tau_k - \tau_{k'}, f)]_{k,k'} \right| \end{aligned} \quad (19)$$

The last equality is true since σ_s^2 and w_j are non-negative. Along with (16), this result implies that among all possible delays τ , it is the delay between two sources that maximizes the contribution of the outer product of their steering vectors to $\bar{\mathbf{R}}(\tau, f)$. This observation is at the core of our method. To better understand its implications, consider the following special case.

C. Special Case: White Source Signal

In this subsection the source signal is assumed to be white, such that $\sigma_s^2(f)$ is constant in f . The weights w_j are all set to 1. Equation (17) can thus be further simplified:

$$\begin{aligned} & \left| [\bar{\mathbf{M}}(\tau, f)]_{k,k'} \right| = \sigma_s^2 \left| \alpha_k \alpha_{k'}^* \sum_{j=0}^{J_f-1} e^{i2\pi\tau j\Delta f (\tau - (\tau_k - \tau_{k'}))} \right| \\ &= \sigma_s^2 \left| \alpha_k \alpha_{k'}^* D_{J_f}(\Delta f (\tau - (\tau_k - \tau_{k'}))) \right| \end{aligned} \quad (20)$$

where:

$$D_n(x) \triangleq \begin{cases} n & x \in \mathbb{Z} \\ \frac{\sin(\pi nx)}{\sin(\pi x)} & x \notin \mathbb{Z} \end{cases} \quad (21)$$

TABLE I
PARAMETERS FOR THE NUMERICAL EXAMPLE

k	τ_k (ms)	α_k
0	0	1
1	2	0.7
2	6	0.5

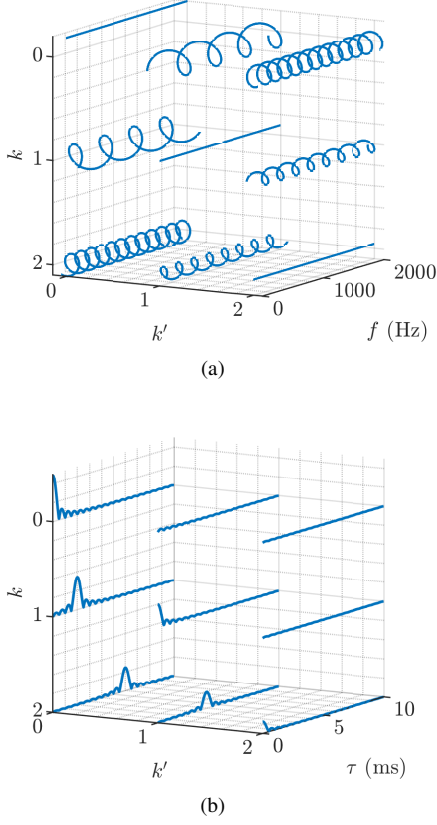


Fig. 1. (a) Entries of the mixing matrix $[\mathbf{M}(f)]_{k,k'}$ for different values of f , with the real and imaginary part of the entries added to the k and k' axes for the purpose of illustrating the complex function. (b) Entries of the transformed mixing matrix $[\mathbf{M}(\tau, 0)]_{k,k'}$ for different values of τ , with the absolute value of the entries added to the k axis for the purpose of illustration.

$D_n(x)$ often arises in Fourier analysis, and is related to the Dirichlet kernel. It has a sinc-like behavior, with a main lobe centered around $x = 0$, and a null-to-null width of $2/n$. Correspondingly, $[\mathbf{M}(\tau, f)]_{k,k'}$ has a main lobe centered at $\tau = \tau_k - \tau_{k'}$, and a width of $2(J_f \Delta f)^{-1}$. Therefore, J_f determines the temporal resolution, which affects the ability to separate reflections with different delays. This result can be used as a guideline for choosing J_f . Note that the same analysis is valid for non white signals, if the weights satisfy $w_j \propto 1/\sigma_s^2(f_j)$.

Next we consider a numerical example where there are $K = 2$ reflections, with parameters summarized in table I. Figure 1a presents the entries of the mixing matrix $\mathbf{M}(f)$. Since the signal is white, the k, k' entry is a complex sinusoidal of the form $\sigma_s^2 \alpha_k \alpha_{k'}^* e^{i2\pi f(\tau_k - \tau_{k'})}$ (see Eqs. (5) and (12)). Its real and imaginary parts are added to the k and k' axes respectively, for the purpose of illustration of the

complex function. This illustration demonstrates that as the delay between two sources is decreased, the period of the corresponding entry, as a function of f , increases. At the extreme, the delay between a source and itself is zero, and so the diagonal entries are constant in frequency.

Figure 1b presents the entries of the transformed mixing matrix $\mathbf{M}(\tau, 0)$, where $\Delta f J_f = 2000\text{Hz}$. In this plot, only the absolute value of the matrix entries is shown, and is added to the k -axis for the purpose of illustration. When τ is equal to a delay between two sources, $\mathbf{M}(\tau, 0)$ is approximately sparse, and the most dominant entry is the one corresponding to the two sources. For other values of τ , all entries of $\mathbf{M}(\tau, 0)$ are relatively small. For $\tau = 0$, $\mathbf{M}(\tau, 0)$ is approximately a diagonal matrix. Since in that case the off diagonal entries can be interpreted as correlations between different sources, this demonstrates the fact that frequency smoothing ($\tau = 0$) performs decorrelation of the sources. If the reflections had the same delay, $\mathbf{M}(0, 0)$ would contain dominant off diagonal entries, and frequency smoothing would have failed to decorrelate the sources. Furthermore, while for $\tau = 0$ all 3 sources are dominant simultaneously, for values of τ that correspond to delays between sources, only a subset of the sources are dominant.

D. Signal-informed Weights Selection

The analysis presented in the previous subsection shows that if the weights $\{w_j\}_j$ are inversely proportional to $\sigma_s^2(f_j)$, the phase alignment transform can effectively separate reflections with different delays. As $\sigma_s^2(f)$ is usually unknown, it must be estimated from the data. A very coarse, yet simple, estimate is given by the trace of $\mathbf{R}(f)$. By neglecting \mathbf{N} in Eq. (13) and substituting Eqs. (12), (5) and (18), we get:

$$\begin{aligned}
 \text{tr}(\mathbf{R}(f)) &\approx \sum_{k,k'} [\mathbf{M}(f)]_{k,k'} \text{tr}(\mathbf{y}(\Omega_k) \mathbf{y}(\Omega_{k'})^H) \\
 &= \sum_{k,k'} \mathbb{E} [s_k(f) s_{k'}^*(f)] \mathbf{y}(\Omega_{k'})^H \mathbf{y}(\Omega_k) \\
 &= \sigma_s^2(f) \sum_{k,k'} e^{-i2\pi f(\tau_k - \tau_{k'})} \alpha_k \alpha_{k'}^* \mathbf{y}(\Omega_{k'})^H \mathbf{y}(\Omega_k) \\
 &= \sigma_s^2(f) \sum_k |\alpha_k|^2 (1 + b_k(f))
 \end{aligned} \tag{22}$$

where:

$$b_k(f) \triangleq 2\Re \left(\sum_{k' > k} \frac{\alpha_{k'}}{\alpha_k} e^{i2\pi f(\tau_k - \tau_{k'})} \mathbf{y}(\Omega_k)^H \mathbf{y}(\Omega_{k'}) \right) \tag{23}$$

and \Re returns the real part of a complex scalar. We argue that b_k is typically small in comparison to 1, since usually the amplitudes decay rapidly. Furthermore, when two reflections have similar amplitudes, it is usually the case that they have very different DOAs, which implies that the inner product of their steering vectors is small [15].

The weights could have another role. Eq. (20) suggests that even if the weights are inversely proportional to σ_s^2 , reflections with delay other than τ may still be dominant in $\mathbf{M}(\tau, f)$, as D_n have strong side lobes. The side lobe levels can be reduced by introducing a window function, at the expense of increasing the width of the main lobe [17].

E. Rank 1 Approximation of Phase Aligned SCM

In the following, the dependence on the frequency f is omitted for brevity. It is important to note that there is no direct access to $\bar{\mathbf{M}}(\tau)$; it is only indirectly observed through $\bar{\mathbf{R}}(\tau)$ as given by Eq. (16). When the proposed transformation succeeds in enhancing a single entry in $\bar{\mathbf{M}}(\tau)$, $\bar{\mathbf{R}}(\tau)$ will be dominated by a single outer product of steering vectors $\mathbf{y}(\Omega)\mathbf{y}(\Omega')^H$. As an outer product is a rank 1 matrix, it is expected that the rank 1 approximation of $\bar{\mathbf{R}}(\tau)$ would perform denoising, i.e. reduce the contribution of $\bar{\mathbf{N}}(\tau)$. The optimal rank 1 approximation (in the least squares sense) of $\bar{\mathbf{R}}(\tau)$ is denoted by $\bar{\mathbf{R}}_1(\tau)$, and is given by truncating its singular value decomposition (SVD):

$$\bar{\mathbf{R}}_1(\tau) = \sigma_\tau \mathbf{u}_\tau \mathbf{v}_\tau^H$$

where σ_τ denotes the first (largest) singular value of $\bar{\mathbf{R}}(\tau)$, and \mathbf{u}_τ and \mathbf{v}_τ denote corresponding left and right singular vectors, respectively. Beside performing denoising, the SVD also separates the two steering vectors, as \mathbf{u}_τ and \mathbf{v}_τ are approximately equal (up to phase) to $\mathbf{y}(\Omega)$ and $\mathbf{y}(\Omega')$, respectively. If there are several reflections with the same delay τ , $\bar{\mathbf{R}}(\tau)$ is still approximately of rank 1 since the dominant entries in $\bar{\mathbf{M}}(\tau)$ all appear at the same column. However, in that case \mathbf{u}_τ is not a single steering vector, but rather a linear combination of the steering vectors of the reflections with delay τ .

V. ALGORITHM DESCRIPTION

This section describes the PHALCOR algorithm for estimating the DOAs and delays of the early reflections. The algorithm is based on the analysis of the first singular vectors of the phase aligned SCM $\bar{\mathbf{R}}(\tau, f)$ that was presented in section IV. The analysis in section IV requires the plane wave decomposition signals to be in the frequency domain. In practice, these are approximated using the short time Fourier transform (STFT) which enables localized analysis in both time and frequency. It is assumed that the window length of the STFT is sufficiently larger than τ_K , such that the multiplicative transfer function (MTF) approximation in the STFT is applicable for Eq. (5) [18]. Note that in the following, τ is the parameter of the phase alignment transform as in Eq. (15), and should not be confused with the time index of the STFT.

The algorithm is performed in two parts. In the first and main part, local time frequency estimates of reflection delays and DOAs are computed; this is performed separately on different regions in the time frequency domain. In the second part, cluster analysis analysis is performed on the local estimates to obtain global estimates that are more robust and accurate.

A. Part 1: Local Time Frequency Estimations

The first, and main, part of the algorithm consists of three steps as described below.

1) *Phase Alignment Transform*: For each time-frequency bin, \mathbf{R} is estimated by replacing the expectation in Eq. (9) with averaging across J_t adjacent bins in time. Then, $\bar{\mathbf{R}}(\tau, f)$ is calculated for $\tau = 0, \Delta\tau, \dots, (J_\tau - 1)\Delta\tau$ using Eq. (15). $\Delta\tau$ dictates the delay estimation resolution, while J_τ dictates the

maximal detectable delay of a reflection. The selection of these parameters is discussed in section V-C, as well as a method to efficiently calculate $\bar{\mathbf{R}}$ using the fast Fourier transform (FFT). The weights are set using the method described in section IV-D:

$$w_j = \frac{W_j}{\text{tr}(\mathbf{R}(f_j))} \quad (24)$$

where W_j is the j 'th sample of a Kaiser window of length J_f , with the β parameter set to 3 [17].

2) *Delay Detection*: The next step is to detect values of τ that are equal to a reflection's delay. Based on the analysis presented in section IV, we suggest the detection of such values of τ by thresholding the following signal:

$$\rho(\tau) = \max_{\Omega' \in \mathbb{S}^2} |\mathbf{y}(\Omega')^H \mathbf{v}_\tau| \quad (25)$$

where \mathbf{v}_τ is a first right singular vector of $\bar{\mathbf{R}}(\tau)$. By the Cauchy-Schwartz inequality, since both \mathbf{v}_τ and $\mathbf{y}(\Omega)$ are unit vectors, $\rho(\tau) \in [0, 1]$ and is equal to 1 if and only if \mathbf{v}_τ is equal (up to phase) to a steering vector. The threshold is set empirically to 0.9.

We denote by $\hat{\Omega}'(\tau)$ the direction that attains the maximum in Eq. (25):

$$\hat{\Omega}'(\tau) = \arg \max_{\Omega' \in \mathbb{S}^2} |\mathbf{y}(\Omega')^H \mathbf{v}_\tau| \quad (26)$$

When τ is equal to a reflection's delay, $\hat{\Omega}'(\tau)$ is an estimate of Ω_0 (the DOA of the direct sound). Note that when τ is equal to a delay between two reflections (and not a delay between a reflection and the direct sound), $\rho(\tau)$ may be high as well, leading to false detections. However, such detections are distinguishable from valid ones, as $\hat{\Omega}'(\tau)$ will be different from Ω_0 ; as detailed in section V-B, the first step of part 2 discards such detections.

3) *DOA Estimation*: The next step is estimating the DOAs of the reflections. This step is performed separately for every τ selected in the previous step. Let \mathbf{u}_τ denote a first left singular vector of $\bar{\mathbf{R}}(\tau)$. According to the analysis presented in section IV, \mathbf{u}_τ is approximately equal to a linear combination of the steering vectors of the reflections with a delay close to τ . If there is only a single such reflection, then its DOA can be estimated using a similar method to that of the direct sound (Eq. (26)). However, in practice there might be several reflections with similar delays. Since their number is expected to be quite small, we utilize sparse recovery to estimate the DOAs. Specifically, we aim to solve the following problem: Find the smallest set of directions $\hat{\Omega}_1, \dots, \hat{\Omega}_S$ and coefficients x_1, \dots, x_S , such that

$$\left\| \sum_{s=1}^S x_s \mathbf{y}(\hat{\Omega}_s) - \mathbf{u}_\tau \right\|^2 \leq \epsilon_u \quad (27)$$

where $\epsilon_u \in (0, 1)$ is a predefined threshold, set experimentally to 0.4. In the context of sparse recovery, the set of vectors $\{\mathbf{y}(\Omega) : \Omega \in \mathbb{S}^2\}$ is known as the dictionary, and its elements are known as atoms. The optimization problem is computationally intractable and cannot be exactly solved in practice. Nevertheless, there has been extensive research on algorithms

that find approximate solutions. In this paper we chose to apply the orthogonal matching pursuit (OMP) algorithm [19]. Although there are more sophisticated sparse recovery algorithms, we chose OMP for several reasons. First, it is simple, and has a low computational cost. Second, although originally designed for finite dictionaries, it is easily extended for our infinite dictionary case. Finally, it is especially suited for our problem by the following argument. Early reflections with similar delays usually have very different DOAs as they typically originate from different walls. If the angle between the DOAs is larger than π/N , the corresponding steering vectors are approximately orthogonal [15] and the projection step in OMP only removes the contribution of steering vectors of DOAs that have already been found, without affecting the rest.

The OMP algorithm is applied on \mathbf{u}_τ for every detected τ . Values of τ where the resulting S is larger than S_{\max} are discarded. The value of S_{\max} was set to 3.

B. Part 2: Cluster Analysis

The input for this part is a list of the local estimates obtained from part 1. Each estimate is a triplet of the form $(\hat{\tau}, \hat{\Omega}, \hat{\Omega}')$, corresponding to the delay of a reflection, its DOA, and the DOA of the direct sound. The goals of this part are: first, to remove outliers, and second, to obtain global estimates for the DOAs and delay of the early reflections.

The first step is to discard estimates where the angle between $\hat{\Omega}'$ and Ω_0 , is larger than some predefined threshold, set empirically to 10 degrees. In general Ω_0 is not known; however, it can be estimated by selecting the peak in the histogram of $\hat{\Omega}'$.

Next, a clustering algorithm is executed on the remaining estimates. We chose the DBSCAN algorithm [20], as it does not require an estimation of the number of clusters, and can automatically detect outliers by not assigning them to any cluster. DBSCAN has two positive parameters ϵ and MINPTS, and operates as follows. Two points are defined as neighbors if the distance between them is less than ϵ . A core point is defined as a point that has MINPTS or more neighbors. A noise point is a non core point, for which none of its neighbors are core points. The algorithm iterates over all the points in the dataset, and assigns two points to the same cluster if one of them is a core point. Noise points are not assigned to any cluster.

The metric we used is the following:

$$d((\tau_a, \Omega_a), (\tau_b, \Omega_b)) = \sqrt{\left(\frac{\angle(\Omega_a, \Omega_b)}{\gamma_\Omega}\right)^2 + \left(\frac{\tau_a - \tau_b}{\gamma_\tau}\right)^2} \quad (28)$$

where γ_Ω and γ_τ are normalization constants, set to 15 degrees and 500 microseconds, respectively. As the metric is normalized, the parameter ϵ is simply set to 1. MINPTS is set to 10 percent of the number of neighbors of the point that has the largest number of neighbors.

Finally, the global estimates are calculated for each cluster by averaging the local estimates of the points assigned to it. The fact that each DOA estimate has an associated delay, is

a major advantage of our method, as it enables the separation of clusters even if they have similar DOAs.

C. Practical Considerations

Avoiding Redundant Processing: The information captured in $\bar{\mathbf{R}}$ contains contributions from a rectangular region in the time-frequency domain of size $J_t \times J_f$. Therefore, it is expected that the results of part 1 of the algorithm would be similar when applied to regions with a large overlap. To reduce computation time, the regions for part 1 are selected with an overlap of 87.5% in frequency.

Acceleration Using the FFT: Part 1 of the algorithm requires the calculation of $\bar{\mathbf{R}}(\tau)$ for a grid of values of τ . Note that if $(\Delta\tau \cdot \Delta f)^{-1} \in \mathbb{N}$, then the sequence

$$(\bar{\mathbf{R}}(j\Delta\tau))_{j=0}^{J_\tau-1} \quad (29)$$

is equal (up to scaling and appropriate zero-padding) to the first J_τ terms of the inverse discrete Fourier transform (taken entry wise) of the sequence

$$(w_j \mathbf{R}(f + j\Delta f))_{j=0}^{J_f-1} \quad (30)$$

so $\bar{\mathbf{R}}$ can be calculated efficiently for the grid of delays using the FFT. A further reduction in the computation time can be achieved using the following identity, obtained directly from Eq. (15) and from the fact that \mathbf{R} is Hermitian:

$$\bar{\mathbf{R}}(\tau) = \bar{\mathbf{R}}(\Delta f^{-1} - \tau)^H \quad (31)$$

Thus, it is sufficient to perform the FFT on only the upper-triangular entries of \mathbf{R} .

Periodicity of the Phase Aligned SCM: It is apparent from Eq. (15) that $\bar{\mathbf{R}}(\tau)$ is periodic with respect to τ , with period Δf^{-1} . This periodicity does not introduce ambiguity in the delay estimation by the following reason. When the STFT window size is T , the frequency resolution of the STFT Δf satisfies $\Delta f \leq 1/T$. Therefore, reflections with delay τ larger than Δf^{-1} necessarily do not satisfy the MTF approximation criteria, since $\tau > T$. This analysis also shows that to avoid unnecessary calculations, J_τ should be chosen such that $J_\tau \Delta\tau < \Delta f^{-1}/2$.

Selecting the Parameters of the Phase Alignment Transform: The calculation of $\bar{\mathbf{R}}$ requires the selection of three parameters: J_f , $\Delta\tau$ and J_τ (see section V-A1). The number of frequency bins J_f , should be chosen to be high enough such that the temporal resolution (given by $(\Delta f J_f)^{-1}$, see section IV-C) is sufficient. For example, if $J_f \Delta f = 2000\text{Hz}$, then the phase alignment transform can separate two reflections if their delays are spaced by more than 0.5ms. However, J_f cannot be set arbitrarily high. First, the frequency independence of steering vectors (see section II-B) is in practice limited to a given band, depending on the geometry of the microphone array. Second, some of our model assumptions may only be valid for bands of limited width. For example, the linear phase assumption in Eq. (5) may, in practice, only hold within a local frequency region.

Once J_f has been determined, a convenient way to set $\Delta\tau$, the delay estimation resolution, is:

$$\Delta\tau = \frac{1}{M\Delta f} \quad (32)$$

where M is an integer that satisfies $M \geq J_f$. This choice guarantees that $\Delta\tau \leq (J_f\Delta f)^{-1}$, and also that $(\Delta\tau \cdot \Delta f)^{-1} \in \mathbb{N}$, so the FFT can be used to calculate $\bar{\mathbf{R}}$. Increasing M beyond J_f would increase the resolution of delay estimation, however it would also increase the computation time of the algorithm.

Finally, J_τ , the number of grid points over τ , should be chosen such that $(J_\tau - 1)\Delta\tau$, the maximal detectable delay, is sufficiently small relative to T , the window length of the STFT, such that the MTF criteria for Eq. (5) holds. From our experience, $(J_\tau - 1)\Delta\tau \approx T/10$ is sufficient.

Maximizing Over the Sphere: Both Eq. (25) and the OMP algorithm require maximizing functions of the form $f(\Omega) = |\mathbf{y}(\Omega)^H \mathbf{x}|$ over the sphere. Note that this is equivalent to finding the maximum of a signal on the sphere whose SFT is given by \mathbf{x} . We use Newton's method to perform this maximization, with initialization obtained by sampling the sphere with a nearly uniform grid of 900 directions [21].

D. Relation To Other Methods

In this section we discuss some similarities between PHALCOR and other methods in array signal processing.

MUSIC and Frequency Smoothing: When $\tau = 0$ and $w_j = 1$, $\bar{\mathbf{R}}(\tau)$ is a frequency-smoothed SCM (as used for example in [22]). Frequency smoothing is a common procedure in source localization in the presence of reverberation, as it can decorrelate signals, which is necessary for subspace methods such as MUSIC. Furthermore, $\bar{\mathbf{R}}(0)$ is positive semi-definite, and therefore its singular vectors are also eigenvectors, so $\hat{\Omega}'(0)$ is the estimate obtained by MUSIC if the signal subspace dimension is set to 1, and $\rho(0)$ is equivalent to the MUSIC spectrum magnitude at this direction. While the frequency smoothing goal is to decorrelate the sources, PHALCOR actually utilizes this correlation to enhance specific reflections.

L1-SVD: Another well known method for source localization that can address correlated sources is L1-SVD [23]. It is based on the observation that the first eigenvectors of the SCM are linear combinations of steering vectors. The DOAs are estimated by decomposing the eigenvectors of the SCM into a sparse combination of steering vectors. This is similar to our method, which decomposes a first left singular vector of the phase aligned SCM to a sparse linear combination of steering vectors. In general, the performance of sparse recovery methods improves as the vectors are more sparse. While in L1-SVD the sparsity is determined by the total number of reflections, in PHALCOR the sparsity is determined by the number of sources at a specific delay, which is significantly lower. Furthermore, like MUSIC, in L1-SVD the number of detectable sources is limited by the number of input channels $((N + 1)^2$ in our case). In PHALCOR, on the other hand, it is possible to detect more reflections than input channels, as each delay is processed independently.

Generalized Cross Correlation: The relations of PHALCOR to MUSIC and L1-SVD is related only to DOA estimation; however PHALCOR is also related to delay estimation methods that are based on generalized cross correlation analysis [24]. It can be shown that the entries of $\bar{\mathbf{R}}(\tau)$ contain a generalized cross correlation between each pair of input channels, at lag τ . Although similar, there are some important distinctions between the two methods. While cross correlation analysis is typically used to estimate the delay between two signals that are observed directly, PHALCOR aims to estimate the delay between multiple signals that are observed indirectly - each input channel is a linear combination of the delayed signals, given by the unknown steering matrix, which is estimated as well.

VI. SIMULATION STUDY

In this section, a simulation study is presented, demonstrating the performance of PHALCOR. First, a detailed analysis of the different steps of the algorithm is presented on a specific test case. Next, a Monte Carlo analysis is presented, demonstrating the robustness of PHALCOR.

A. Simulation Setup

The setup of the simulations, common to both the case study and the Monte Carlo study, is as follows. An acoustic scene that consists of a speaker and a rigid spherical microphone array in a shoe box room, was simulated using the image method [16]. The speech signal is a 5 seconds sample, drawn randomly from the TSP Speech Database [25]. The array has 32 microphones, and a radius of 4.2 cm (similar to the Eigenmike [26]), facilitating plane wave decomposition with spherical harmonics order $N = 4$. The microphone signals were sampled at 48 KHz. Sensor noise was added, such that the direct sound to noise ratio is 30dB. The positions of the speaker and the array were chosen at random inside the room, with the restriction that the distance between each other, and to the room boundaries is no less than 1 meter. Three different rooms sizes are considered. Their dimensions and several acoustic parameters, are presented in table II.

B. Methodology

The signals recorded by the microphones were used to compute $\mathbf{a}_{\text{nm}}(f)$ as detailed in section II-B. An STFT was applied to the PWAD coefficients signals using a Hanning window of 8192 samples, and an overlap of 75%. A frequency range of [500, 5000] Hz was chosen for the analysis. The number of time bins used for averaging, J_t , was set to 6, while the number of frequency bins used for the phase alignment transform, J_f , was set such that $J_f\Delta f = 2000$ Hz. The delay resolution $\Delta\tau$ was set to 83.33 microseconds (equivalent to setting $M = 2048$ in Eq. (32)), while J_τ was chosen such that the maximal delay is 20 ms. With these parameters, PHALCOR, detailed in section V, was applied to the simulated data.

The MUSIC algorithm [22] was applied as a baseline reference method for DOA estimation, by selecting the peaks

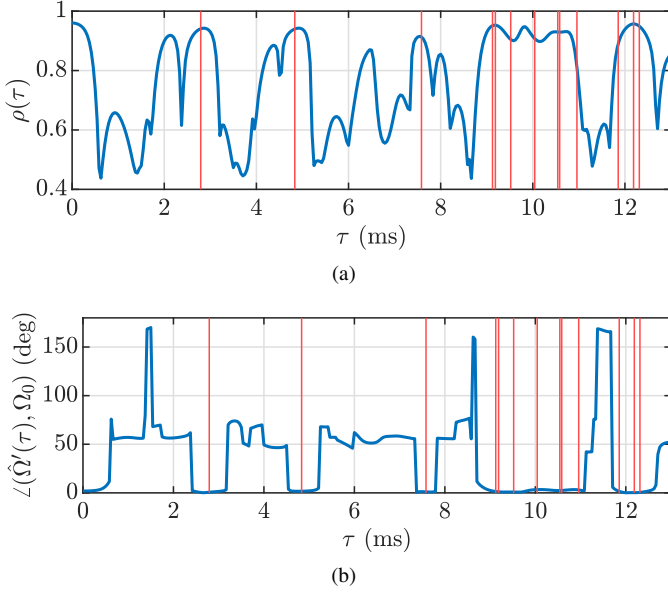


Fig. 2. (a) ρ as function of τ (Eq. (25)). (b) The angle between $\hat{\Omega}'(\tau)$ (Eq. (26)) and true direction of the direct sound, as a function of τ . The red vertical lines correspond to the true delays of the reflections.

in the MUSIC spectrum $\|\mathbf{y}(\Omega)^H \mathbf{U}\|$, where \mathbf{U} is matrix whose columns are orthonormal eigenvectors that correspond to the L largest eigenvalues of the time and frequency smoothed SCM. The time and frequency smoothing parameters are the same as in PHALCOR. The dimension of the signal subspace L was determined using the SORTe method [27]. To reduce false positives, peaks for which the MUSIC magnitude height is lower than 0.75 are discarded. The local estimates are clustered using DBSCAN, to obtain global estimates. The delays of the detected reflections are estimated using the following method, which is similar to the one proposed in [5]. First, each reflection signal is estimated by solving Eq. (6) for \mathbf{s} in the least squares sense. Then, the delay of the k 'th reflections is estimated by selecting the maximum of the cross correlation values between s_k and s_0 . Note that in contrast to PHALCOR, the delays are estimated only after the clustering.

For both PHALCOR and the baseline method, we consider a detected reflection as a true positive if its delay and DOA match simultaneously the delay and DOA of a true reflection. The matching tolerance was chosen to be 500 μs for the delay, and 15 degrees for the DOA. The probability of detection (PD) at a given maximal delay t is defined as the fraction of true positive detections with delay less than t , out of the total expected number of reflections with delay smaller than t . The false positive rate (FPR) at a given maximal delay t is defined as the fraction of false positive detections with delay less than t , out of the total number of detections with delay smaller than t .

C. Results of a Case Study

The test case presented in this section is of a female speaker in the medium sized room. There are $K = 31$ reflections with a delay less than 20 ms in this case.

Figures 2a and 2b illustrate the delay detection routine, as detailed in section V-A2. Figure 2a shows $\rho(\tau)$ as a function of

τ . Since $\rho(\tau)$ measures the similarity between \mathbf{v}_τ , a first right singular vector of $\hat{\mathbf{R}}(\tau)$, and a steering vector (see Eq. (25)), it is high for values of τ that are close to a reflection's delay, indicated on the plot using red vertical lines. There are also values of τ that are not close to a reflection's delay, for which $\rho(\tau)$ is high. These correspond to delays between two reflections (as opposed to delays between a reflection and the direct sound). For example, the peak near $\tau = 2\text{ms}$, corresponds to the delay between the second and third reflections, whose delays are about 3ms and 5ms, respectively. As discussed in section V-A2, such cases may be identified by testing $\angle(\hat{\Omega}'(\tau), \Omega_0)$, the angle between the DOA of the steering vector that is most similar to \mathbf{v}_τ , and the DOA of the direct sound. As shown in figure 2b, $\angle(\hat{\Omega}'(\tau), \Omega_0)$ is small for values of τ which are close to a reflection's delay, and high otherwise. Therefore, false detections such as $\tau = 2\text{ms}$, will be discarded during the first step of the second part of the algorithm.

Figure 3 illustrates the process of DOA estimation, as detailed in section V-A3. Each plot shows a different function on the sphere, which is projected onto the 2D page using the Hammer projection. In the top row, the function $|\mathbf{y}(\Omega)\mathbf{v}_\tau|$ is shown, where each column corresponds to a different value of τ . Recall that when τ equals a reflection's delay, we expect the direction that maximizes the response to be that of the direct sound. Indeed, as τ varies across columns, the location of the peak remains, and is equal to Ω_0 , the DOA of the direct sound.

In the middle row, the function $|\mathbf{y}(\Omega)\mathbf{u}_\tau|$ is shown. It is similar to the top row, except that a first left singular vector is used instead of a right one. Recall that when τ is a reflection's delay, \mathbf{u}_τ is approximately equal to a linear combination of the steering vectors of reflections with delays of approximately τ . When the DOAs are sufficiently separated, they can be identified as peaks in $|\mathbf{y}(\Omega)\mathbf{u}_\tau|$. For τ_1 and τ_2 , only one such peak is apparent, and its location matches the DOA of the corresponding reflection. When $\tau = \tau_4$, it is apparent that there are two dominant peaks, at directions Ω_4 and Ω_5 . This is due to the fact that the 4th and 5th reflections have similar delays. Similarly, since the 8th and 9th reflections have similar delays, when $\tau = \tau_8$ the two peaks correspond to Ω_8 and Ω_9 .

Figure 3 demonstrates the effectiveness of PHALCOR in separating reflections from the direct sound, as well as reflections with different delays. This is in contrast to the MUSIC spectrum (shown the same figure), which shows only a few peaks, which are much less separable; as a result, fewer reflections are potentially identified.

Figures 4 and 5 present the local estimates obtained with PHALCOR (as detailed in section V-A) and the baseline methods (as detailed in section VI-B), respectively. It is apparent that, compared to the baseline, PHALCOR is able to detect significantly more reflections. PHALCOR detected successfully 29 reflections, while the MUSIC based method could only detect 8 (not including the direct sound). Furthermore, figures 4 and 5 illustrate the advantage of simultaneously estimating DOA and delay for cluster analysis.

TABLE II
ROOM PARAMETERS USED IN THE SIMULATION STUDY

Room	Dimensions (m)	Reverberation Time (s)	Critical Distance (m)	Average Distance Between Source and Array (m)	Average Number of Reflections With Delay Smaller Than 20 ms
Small	$5 \times 4 \times 2.5$	0.6	0.5	1.7	52
Medium	$7 \times 5 \times 3$	0.8	0.7	2.5	31
Large	$10 \times 7 \times 3.5$	1.1	1	3.8	21

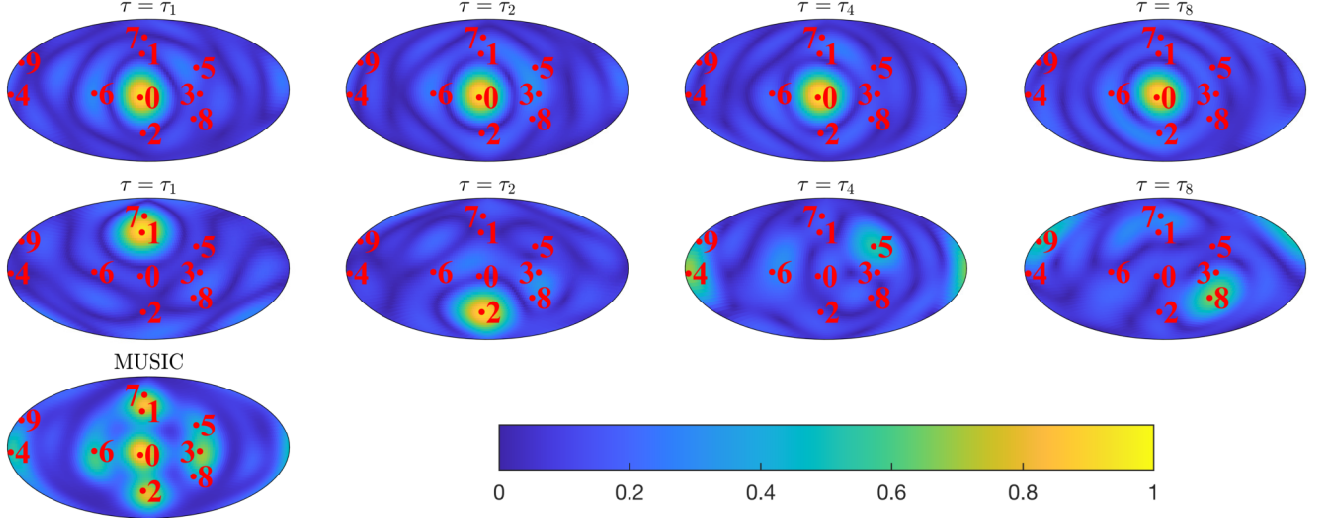


Fig. 3. Top row: $|\mathbf{y}(\Omega)^H \mathbf{v}_\tau|$ as a function of Ω , for different values of τ . Middle row: $|\mathbf{y}(\Omega)^H \mathbf{u}_\tau|$ as a function of Ω , for different values of τ . Bottom figure: MUSIC spectrum. The red markers correspond to ground truth DOAs $\Omega_0, \dots, \Omega_9$, with the numbers indicating the index. The Hammer projection is used to project the sphere onto the plane.

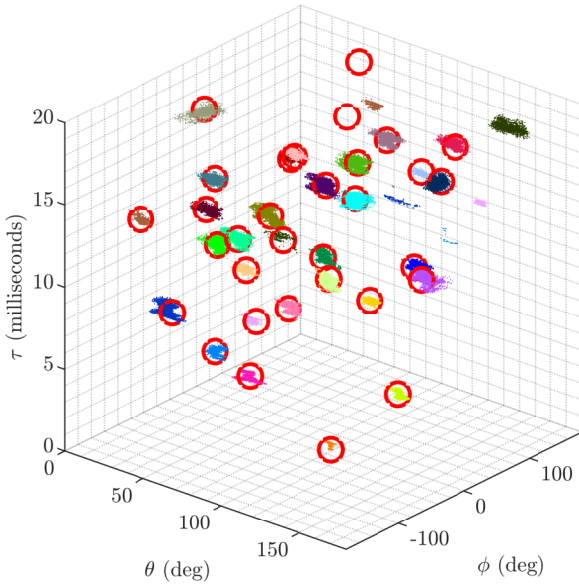


Fig. 4. Local DOA and delay estimates obtained by PHALCOR, colored by assigned cluster number. The θ and ϕ axes are for elevation and azimuth, respectively. The red circles correspond to true DOAs and delays.

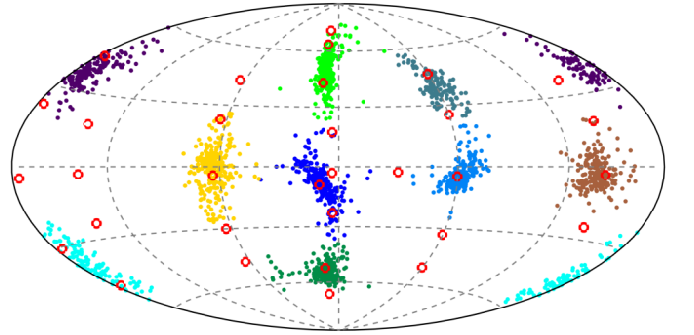


Fig. 5. Local DOA estimates obtained by MUSIC, colored by assigned cluster number. The red circles correspond to $\Omega_0, \dots, \Omega_{31}$. The Hammer projection is used to project the sphere onto the plane.

D. Monte Carlo Analysis

The simulation described above is repeated 50 times for each of the 3 rooms, varying the speakers, their location, and the microphone array location, as detailed in section VI-A. Figure 6 presents PD and FPR, as defined in section VI-B, for different values of t , the maximum delay of the identified reflections. Compared with the baseline method, the performance of PHALCOR is significantly better, both in terms of probability of detection and false positive rates, by a factor ranging from 3 to 20. As the delay of a reflection increases, the probability of detection decreases. This is since

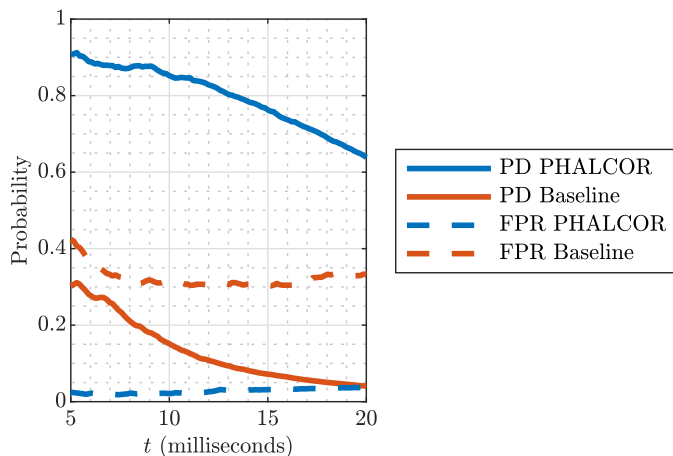


Fig. 6. PD and FPR, as defined in section VI-B, of PHALCOR and the baseline method

TABLE III
AVERAGE ESTIMATION ERRORS FOR THE ENTIRE MONTE CARLO SIMULATION

Method	DOA Error RMS (deg)	Delay Error RMS (μ s)	Average Number Of True Positive Detections
PHALCOR	4.3	77	24.8
Baseline	6.5	43	3.8

later reflections usually have lower amplitudes. Furthermore, the reflection density is higher as the delay increases, making it more difficult to separate the reflections spatially.

The root mean square (RMS) for DOA and delay estimation errors for each method are computed and averaged for all the estimates in this Monte Carlo simulation, and are presented in table III. The RMS is calculated excluding the direct sound. The table shows that the performance in terms of estimation error is comparable between the two methods, but note that the errors are calculated only on true positive detections, which are considerably more frequent in PHALCOR, as is evident from figure 6 and the last column of table III.

VII. CONCLUSIONS

In this work, PHALCOR, a novel method method for estimating the DOA and delay of early reflections, is proposed. The method is based on a phase alignment transform of the spatial correlation matrices, which enables the detection of reflections with similar DOAs. A simulation study showed that the proposed method is able to detect and localize a large number of reflections compared to existing methods. The estimation of reflection amplitudes and the validation of the method performance on measured data are proposed for future research.

REFERENCES

- [1] K. Kowalczyk, S. Kacprzak, and M. Ziółko, "On the extraction of early reflection signals for automatic speech recognition," in *2017 IEEE 2nd International Conference on Signal and Image Processing (ICSIP)*. IEEE, 2017, pp. 351–355.
- [2] Y. Peled and B. Rafaely, "Method for dereverberation and noise reduction using spherical microphone arrays," in *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2010, pp. 113–116.
- [3] E. Vincent, N. Bertin, R. Gribonval, and F. Bimbot, "From blind to guided audio source separation: How models and side information can improve the separation of sound," *IEEE Signal Processing Magazine*, vol. 31, no. 3, pp. 107–115, 2014.
- [4] H. A. Javed, A. H. Moore, and P. A. Naylor, "Spherical microphone array acoustic rake receivers," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 111–115.
- [5] E. Mabande, K. Kowalczyk, H. Sun, and W. Kellermann, "Room geometry inference based on spherical microphone array eigenbeam processing," *The Journal of the Acoustical Society of America*, vol. 134, no. 4, pp. 2773–2789, 2013.
- [6] J. Catic, S. Santurette, and T. Dau, "The role of reverberation-related binaural cues in the externalization of speech," *The Journal of the Acoustical Society of America*, vol. 138, no. 2, pp. 1154–1167, 2015.
- [7] M. Vorländer, *Auralization: fundamentals of acoustics, modelling, simulation, algorithms and acoustic virtual reality*. Springer Science & Business Media, 2007.
- [8] V. Pulkki, S. Delikaris-Manias, and A. Politis, *Parametric time-frequency domain spatial audio*. Wiley Online Library, 2018.
- [9] P. Coleman, A. Franck, P. Jackson, R. J. Hughes, L. Remaggi, F. Melchior *et al.*, "Object-based reverberation for spatial audio," *Journal of the Audio Engineering Society*, vol. 65, no. 1/2, pp. 66–77, 2017.
- [10] H. Kuttruff, *Room acoustics*. Crc Press, 2016.
- [11] H. Sun, E. Mabande, K. Kowalczyk, and W. Kellermann, "Localization of distinct reflections in rooms using spherical microphone array eigenbeam processing," *The Journal of the Acoustical Society of America*, vol. 131, no. 4, pp. 2828–2840, 2012.
- [12] B. Jo and J.-W. Choi, "Robust localization of early reflections in a room using semi real-valued eb-esprit with three recurrence relations and laplacian constraint," in *International Commission for Acoustics (ICA)*. International Commission for Acoustics (ICA), 2019.
- [13] P. K. T. Wu, N. Epain, and C. Jin, "A dereverberation algorithm for spherical microphone arrays using compressed sensing techniques," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012, pp. 4053–4056.
- [14] T. Shlomo and B. Rafaely, "Blind localization of early room reflections from reverberant speech using phase aligned spatial correlation," in *INTERSPEECH 2020, submitted for publication*.
- [15] B. Rafaely, *Fundamentals of spherical array processing*. Springer, 2015, vol. 8.
- [16] J. B. Allen and D. A. 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.
- [17] A. V. Oppenheim, *Discrete-time signal processing*. Pearson Education India, 1999.
- [18] Y. Avargel and I. Cohen, "On multiplicative transfer function approximation in the short-time fourier transform domain," *IEEE Signal Processing Letters*, vol. 14, no. 5, pp. 337–340, 2007.
- [19] T. T. Cai and L. Wang, "Orthogonal matching pursuit for sparse signal recovery with noise," *IEEE Transactions on Information theory*, vol. 57, no. 7, pp. 4680–4688, 2011.
- [20] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Kdd*, vol. 96, no. 34, 1996, pp. 226–231.
- [21] J. Fliege and U. Maier, "A two-stage approach for computing cubature formulae for the sphere," in *Mathematik 139T, Universität Dortmund, Fachbereich Mathematik, Universität Dortmund, 44221*. Citeseer, 1996.
- [22] D. Khaykin and B. Rafaely, "Coherent signals direction-of-arrival estimation using a spherical microphone array: Frequency smoothing approach," in *2009 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*. IEEE, 2009, pp. 221–224.
- [23] D. Malioutov, M. Cetin, and A. S. Willsky, "A sparse signal reconstruction perspective for source localization with sensor arrays," *IEEE transactions on signal processing*, vol. 53, no. 8, pp. 3010–3022, 2005.
- [24] J. Hassab and R. Boucher, "Optimum estimation of time delay by a generalized correlator," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 27, no. 4, pp. 373–380, 1979.
- [25] P. Kabal, "Tsp speech database," *McGill University, Database Version*, vol. 1, no. 0, pp. 09–02, 2002.
- [26] M. Acoustics, "Em32 eigenmike microphone array release notes (v17.0)," 25 Summit Ave, Summit, NJ 07901, USA, 2013.
- [27] K. Han and A. Nehorai, "Improved source number detection and direction estimation with nested arrays and ulas using jackknifing," *IEEE Transactions on Signal Processing*, vol. 61, no. 23, pp. 6118–6128, 2013.