



VIRTUAL NAVIGATION VIA HIGHER ORDER DISTRIBUTED SOUND SOURCES

Thushara D. Abhayapala^{1*} Jihui (Aimee) Zhang^{2,1} Shaoheng Xu¹
 David Lou Alon³ Zamir Ben-Hur³ Prasanga N. Samarasinghe¹

¹ The Australian National University, Australia

² University of Southampton, UK

³ Reality Labs Research at Meta, USA

ABSTRACT

With the rise of virtual reality, there is a demand for recording and recreating real-world experiences that allow users to move throughout audio-visual scenes. Sound field translation achieves this by building an equivalent environment of virtual sources to recreate the recording spatially. However, combining multiple microphone array recordings and dealing with mix fields of exterior and interior together is difficult. In this paper, we explain a novel method of virtual navigation by representing complex sound fields by distributed virtual higher order sound sources. The technique combine recordings from spatially separated microphones and decompose them to a sparse representation of multiple higher order sources. The method can be thought of sound field decomposition into a grid of higher order sources. We show the method can be used in VR applications involving navigation within complex sound scenes.

Keywords: *higher order sound sources, virtual navigation, virtual sources*

1. INTRODUCTION

As the concept of metaverse and the applications of virtual reality gains popularity, there is a growing need to capture

*Corresponding author: u9701943@anu.edu.au.

Copyright: ©2023 Thushara D. Abhayapala et al. This is an open-access article distributed under the terms of the Creative Commons Attribution 3.0 Unported License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

and replicate real-life experiences, enabling users to navigate through audio-visual environments [1]. In this paper, we are interested in recreating spatial audio environment for virtual navigation. Whilst there have been recent advances in this area [2–4], combining multiple microphone array recordings and dealing with mix fields of exterior and interior together have not been treated well.

Recording a three dimensional sound scene with multiple active sound sources in a reverberant room, for the purpose of accurately recreating, needs a dense sampling of the space by a large number of distributed microphones or microphone arrays. It is impractical to place a large number of dedicated microphones to satisfy the sampling theorem over space. However, recent advances in smart devices and wearable devices such as smart glasses with microphones can be potentially used to record and combine spatial sound over the space.

Typically, sound field navigation in VR setting is implemented by finding an equivalent representation of the recorded sound field. In [2], a mix of near-field and far-field virtual sources exterior to the recording area with sparse constrained expansion is used to binaurally reproduce the soundfield enabling translation. This method has been perceptually validated in [4] and extended with bilateral-ambisonic reproduction [5]. The authors of [6], proposed a method to reconstruct room impulse responses over extended domains for navigable sound field reproduction. In [7], an object based six-degrees-of-freedom rendering of sound scenes is proposed using the recorded signals from multiple spherical microphone arrays. Another sparse method is proposed in [8] where compressive sensing and image-source models were used. A paramet-

ric decomposition is used in [3] to translate higher order spherical harmonic sound scene recording. In [9], two spherical microphone arrays have been used to capture sound and combine for interpolation.

None of the above work describe a method to combine multiple microphone arrays with different orders together nor deal with mixed sound field of exterior and interior field together. In this paper, our goal is to find an equivalent representation that will enable to (i) combine recordings from spatially distributed microphones and microphone arrays including smart device embedded microphones, and (ii) navigate around a spatial sound scene with mix exterior and interior recording devices/sound sources.

2. PROBLEM FORMULATION

Let there are Q microphone arrays/devices with the q^{th} array of order N_q located at position \mathbf{x}_q , $q = 1, \dots, Q$. Assume that the received spatial signal by the q^{th} microphone array is recorded in the spherical harmonic domain as coefficients $\beta_{nm}^{(q)}(k)$, where $n = 0, \dots, N_q$, $m = -n, \dots, n$, and k is the wave number. The harmonic coefficient vector at \mathbf{x}_q is defined as $\beta^{(q)}(k) = [\beta_{00}^{(q)}(k), \dots, \beta_{N_q N_q}^{(q)}(k)]$. By stacking all Q received spherical harmonic coefficient vectors, we define the total recorded coefficient vector as

$$\beta(k) = [\beta^{(1)}(k), \beta^{(2)}(k), \dots, \beta^{(Q)}(k)]^T, \quad (1)$$

where $(\cdot)^T$ represents the transpose operator. The problem addressed in this paper is, given $\beta(k)$, find an equivalent sparse representation to reconstruct the recorded acoustic scene that enable binaural spatial sound reproduction.

3. HIGHER ORDER SOURCES

In this section we outline the notion of *virtual higher order sources* guided by the concept of higher order sources (HOS) introduced in [10, 11].

An arbitrary higher order sound source located at an origin will produce a sound field $S(\mathbf{y}, k)$ at a location \mathbf{y} from the origin

$$S(\mathbf{y}, k) = \sum_{n=0}^N \sum_{m=-n}^n \alpha_{nm}(k) h_n(k|\mathbf{y}|) Y_{nm}(\hat{\mathbf{y}}), \quad (2)$$

where $\hat{\mathbf{y}}$ is a unit vector in the direction of \mathbf{y} , $h_n(\cdot)$ are the spherical Hankel functions of the first kind of order n ,

$Y_{nm}(\cdot)$ are the spherical harmonics of order n and mode m , N is the order of the source and $\alpha_{nm}(k)$ are the directivity coefficients in the spherical harmonic domain. In the case of a point source, i.e., a zeroth order source, (2) is reduced to

$$\begin{aligned} S(\mathbf{y}, k) &= \frac{e^{ik|\mathbf{y}|}}{4\pi|\mathbf{y}|} = \frac{ik}{4\pi} h_0(k|\mathbf{y}|) \\ &= \frac{ik}{\sqrt{4\pi}} h_0(k|\mathbf{y}|) Y_{00}(\hat{\mathbf{y}}). \end{aligned} \quad (3)$$

Thus, the sound field produce by an N^{th} order sound source at \mathbf{y} is proportional to

$$S(\mathbf{y}, k) \propto h_n(k|\mathbf{y}|) Y_{nm}(\hat{\mathbf{y}}). \quad (4)$$

It has argued in [11] that the above form is less practical due to numerical issues when kr tends to zero due to the Hankel function. Thus, a more realizable higher order loudspeaker is derived by considering a physical size of radius a for the source with vibrating surface velocity. Thus, in this work we use numerically stable, physically supported higher order source of the form

$$S(\mathbf{y}, k) \propto \frac{i\rho c}{k} \frac{h_n(k|\mathbf{y}|)}{h'_n(ka)} Y_{nm}(\hat{\mathbf{y}}), \quad (5)$$

where ρ is the density of air, c is the speed of sound and $h'_n(\cdot)$ represent the derivative of the spherical Hankel function with respect to its argument.

Now, we can write the sound field at a location \mathbf{y} due to a higher order source at a location \mathbf{d} with respect to an origin as

$$S(\mathbf{y}, k) \propto \frac{i\rho c}{k} \frac{h_n(k|\mathbf{y} - \mathbf{d}|)}{h'_n(ka)} Y_{nm}\left(\frac{\mathbf{y} - \mathbf{d}}{|\mathbf{y} - \mathbf{d}|}\right). \quad (6)$$

4. EQUIVALENT GRID OF VIRTUAL HOS

Consider a suitable grid of L *virtual higher order sources* with the ℓ^{th} , $\ell = 1, \dots, L$ source located at \mathbf{d}_ℓ . Then, the received signal at a location $\mathbf{x}^{(q)}$ with respect to the origin of the q^{th} microphone array is

$$\begin{aligned} S_{\text{eq}}(\mathbf{x}^{(q)}, k) &= \sum_{\ell=1}^L \sum_{n=0}^N \sum_{m=-n}^n \frac{h_n(k|\mathbf{y}_q + \mathbf{x}^{(q)} - \mathbf{d}_\ell|)}{h'_n(ka)} \\ &\times Y_{nm}\left(\frac{\mathbf{y}_q + \mathbf{x}^{(q)} - \mathbf{d}_\ell}{|\mathbf{y}_q + \mathbf{x}^{(q)} - \mathbf{d}_\ell|}\right) w_{nm}^{(\ell)}(k), \end{aligned} \quad (7)$$

where $w_{nm}^{(\ell)}(k)$ are the weights of the ℓ^{th} source of order n and mode m . Note that the factor $i\rho c/k$ in (6) has been

incorporated in to the weight $w_{nm}^{(\ell)}(k)$ in (7). We write (7) in the matrix as

$$S_{\text{eq}}(\mathbf{x}^{(q)}, k) = \mathbf{h}(\mathbf{x}^{(q)}, k) \mathbf{w}(k), \quad (8)$$

where $\mathbf{h}(\mathbf{x}^{(q)}, k)$ is a $1 \times L(N+1)^2$ vector with $(\ell-1)(N+1)^2 + n^2 + n + m + 1$ element is given by

$$\frac{h_n(k|\mathbf{y}_q + \mathbf{x}^{(q)} - \mathbf{d}_\ell|)}{h'_n(ka)} Y_{nm} \left(\frac{\mathbf{y}_q + \mathbf{x}^{(q)} - \mathbf{d}_\ell}{|\mathbf{y}_q + \mathbf{x}^{(q)} - \mathbf{d}_\ell|} \right),$$

and

$$\mathbf{w}(k) = [w_{00}^{(1)}(k), w_{1(-1)}^{(1)}(k), w_{10}^{(1)}(k), \dots, w_{NN}^{(1)}(k), \dots, \dots, w_{00}^{(L)}(k), \dots, w_{NN}^{(L)}(k)]^T.$$

We can also write the sound field at $\mathbf{x}^{(q)}$ using the measured sound field coefficients $\beta^{(q)}(k)$ by the q^{th} microphone array as

$$S_{\text{me}}(\mathbf{x}^{(q)}, k) = \sum_{\nu=0}^{N_q} \sum_{\mu=-\nu}^{\nu} \beta_{\nu\mu}^{(q)}(k) j_\nu(k|\mathbf{x}^{(q)}|) Y_{\nu\mu}(\hat{\mathbf{x}}^{(q)}), \quad (9)$$

where $j_n(\cdot)$ are the spherical Bessel functions. We write (9) in matrix form as

$$S_{\text{me}}(\mathbf{x}^{(q)}, k) = \mathbf{j}(\mathbf{x}^{(q)}, k) \boldsymbol{\beta}^{(q)}(k), \quad (10)$$

where

$$\mathbf{j}(\mathbf{x}^{(q)}, k) = \begin{bmatrix} j_0(k|\mathbf{x}^{(q)}|) Y_{00}(\hat{\mathbf{x}}^{(q)}) \\ \vdots \\ j_{N_q}(k|\mathbf{x}^{(q)}|) Y_{N_q N_q}(\hat{\mathbf{x}}^{(q)}) \end{bmatrix}^T.$$

We can use (7) and (9) to connect the measured coefficients $\beta^{(q)}$ to weights of the virtual higher order loudspeakers. Note that there are multiple way of deriving this expression, including equating the two equation over multiple spheres/points around the q^{th} microphone array origin and performing either pseudo inverse or orthogonal decomposition by multiplying conjugate spherical harmonics ($Y_{nm}^*(\mathbf{x}^{(q)})$) and integration (approximated via a suitable summation). As one of these possibilities, below, we evaluate $S_{\text{eq}}(\mathbf{x}^{(q)}, k)$ and $S_{\text{me}}(\mathbf{x}^{(q)}, k)$ over P_q , $q = 1, \dots, Q$ points for all microphone locations $\mathbf{x}_p^{(q)}$.

By evaluating (8) at $\mathbf{x}_p^{(q)}$ for $p = 1, \dots, P_q$ over $q = 1, \dots, Q$, we get

$$\mathbf{S}_{\text{eq}} = \mathbf{H} \mathbf{w}, \quad (11)$$

where $\mathbf{S}_{\text{eq}} =$

$$[S_{\text{eq}}(\mathbf{x}_1^{(1)}) \dots S_{\text{eq}}(\mathbf{x}_{P_1}^{(1)}) \dots S_{\text{eq}}(\mathbf{x}_1^{(Q)}) \dots S_{\text{eq}}(\mathbf{x}_{P_Q}^{(Q)})]^T,$$

and $(P_1 + P_2 + \dots + P_Q) \times L(N+1)^2$ matrix

$$\mathbf{H} = \begin{bmatrix} \mathbf{h}(\mathbf{x}_1^{(1)}) \\ \vdots \\ \mathbf{h}(\mathbf{x}_{P_q}^{(1)}) \\ \vdots \\ \mathbf{h}(\mathbf{x}_1^{(Q)}) \\ \vdots \\ \mathbf{h}(\mathbf{x}_{P_Q}^{(Q)}) \end{bmatrix}.$$

Similarly, using (10), we obtain

$$\mathbf{S}_{\text{me}} = \begin{bmatrix} \mathbf{j}(\mathbf{x}_1^{(1)}) \boldsymbol{\beta}^{(1)} \\ \vdots \\ \mathbf{j}(\mathbf{x}_{P_1}^{(1)}) \boldsymbol{\beta}^{(1)} \\ \vdots \\ \mathbf{j}(\mathbf{x}_1^{(Q)}) \boldsymbol{\beta}^{(Q)} \\ \vdots \\ \mathbf{j}(\mathbf{x}_{P_1}^{(Q)}) \boldsymbol{\beta}^{(Q)} \end{bmatrix}, \quad (12)$$

where $\mathbf{S}_{\text{me}} =$

$$[S_{\text{me}}(\mathbf{x}_1^{(1)}) \dots S_{\text{me}}(\mathbf{x}_{P_1}^{(1)}) \dots S_{\text{me}}(\mathbf{x}_1^{(Q)}) \dots S_{\text{me}}(\mathbf{x}_{P_Q}^{(Q)})]^T.$$

Note that we omit the frequency dependency k for brevity in the above equations and for the rest of the paper.

4.1 Sparse Solution

The aim is to find an equivalent higher order source weight vector \mathbf{w} such that \mathbf{S}_{eq} in (11) as close as possible to the \mathbf{S}_{me} which is derived from the measurements. We can formulate this problem as a least squares (l_2) optimisation as

$$\min_{\mathbf{w}} \|\mathbf{S}_{\text{me}} - \mathbf{H} \mathbf{w}\|_2^2, \quad (13)$$

or to obtain a sparse solution (i.e., l_1 optimisation) as

$$\min_{\mathbf{w}} \left(\|\mathbf{S}_{\text{me}} - \mathbf{H} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1 \right), \quad (14)$$

where $\|\cdot\|_p$ is the p-norm and λ is a suitable turning parameter.

In our application, we assume that the number of independent sources within a room (a recording) environment is finite and sparse. Also, the number of recording microphones (including higher order) are limited and randomly placed. Therefore, we seek a sparse solution to the above problem using a large number of higher order source dictionary. Having a sparse solution will assist an efficient implementation of any binaural navigation implementation.

4.2 Choice of HOS Grid

There are many possible geometries for the grid of virtual higher order sources. One approach is a grid of circular or spherical shell (including multiple circles) encompassing the room. However, this geometry will not support navigation around actual sound sources in the room as the equivalent field will be an exterior sound field. Another approach is to have a grid of higher order sources distributed over the whole volume of the room. This approach can support navigation around sound sources however will struggle to represent reflections from walls. Thus, in this paper, we use a grid of higher order sources distributed inside the room and a ring/sphere of point sources with sufficiently large radius encompassing the room to represent reverberations. We are only considering the proof of the concept in this paper and leave further investigation for future work.

4.3 Point Source vs HOS

Since any higher order source can be represented by a set of point sources located over a region of space, one can question the benefit of having higher order sources instead of point sources. The reasons are (i) HOS grid is less dense than an equivalent point source grid (e.g., a single N th order source needs at least $(N + 1)^2$ point sources), (ii) HOS grid is computationally efficient and numerically stable to implement in virtual navigation applications, and (iii) HOS are natural to model directional sources. However, the number of parameters (weights) are equivalent in both HOS and point source grids.

4.4 Sparse Algorithms

There are many sparse algorithms present in literature. They are mainly based on matching pursuit [12] or basis pursuit [13] and their variations. For the proposed concept in this paper, we can use any one of them including *least absolute shrinkage and selection operator* (LASSO) [14]

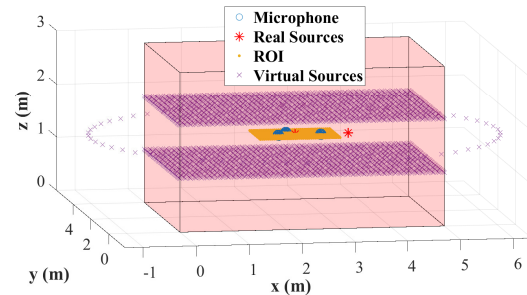


Figure 1: Simulation setup in case 2, where the real sound sources are marked as red *, microphones are marked as blue o, virtual higher order sources are marked as purple x, and the ROI is marked as yellow.

and *iteratively re-weighted least squares* (IRLS) [15]. Following up on our recent work [2], we use IRLS method in the simulation section of this paper. However, we argue that any of the other sparse algorithms can be used instead of IRLS.

5. SIMULATIONS

In this section, we validate the sound navigation using the proposed virtual HOS method, and compare the results using traditional least squares (ℓ_2) solution and IRLS sparse solution.

5.1 Simulation setups

In this numerical simulation, we setup a room with a dimension of $5 * 4 * 3$ m. The reflection coefficients for the four walls are $[0.6, 0.6, 0.6, 0.6]$, and we assume there are no reflections on the ceiling and floor. The region of interest (ROI) is a rectangular area ($1.5 * 1.2$ m) with center at $[2.5, 2, 1.5]$. As shown in Figure 1, three higher order microphone arrays (HOM) with a right 30-60-90 degree triangle have been placed in the middle of ROI to record the sound field. Each HOM is an open spherical array with 32 microphones based on the eigenmike [16] placement. The distance between HOM1 and HOM2 is 0.4 m in case 1, and 0.8 m in case 2. A signal-to-noise ratio of 40 dB white Gaussian noise is added at each microphone recording.

Two sound sources have been simulated, where one source locates in the middle of three HOMs at $[2.5, 2, 1.5]$

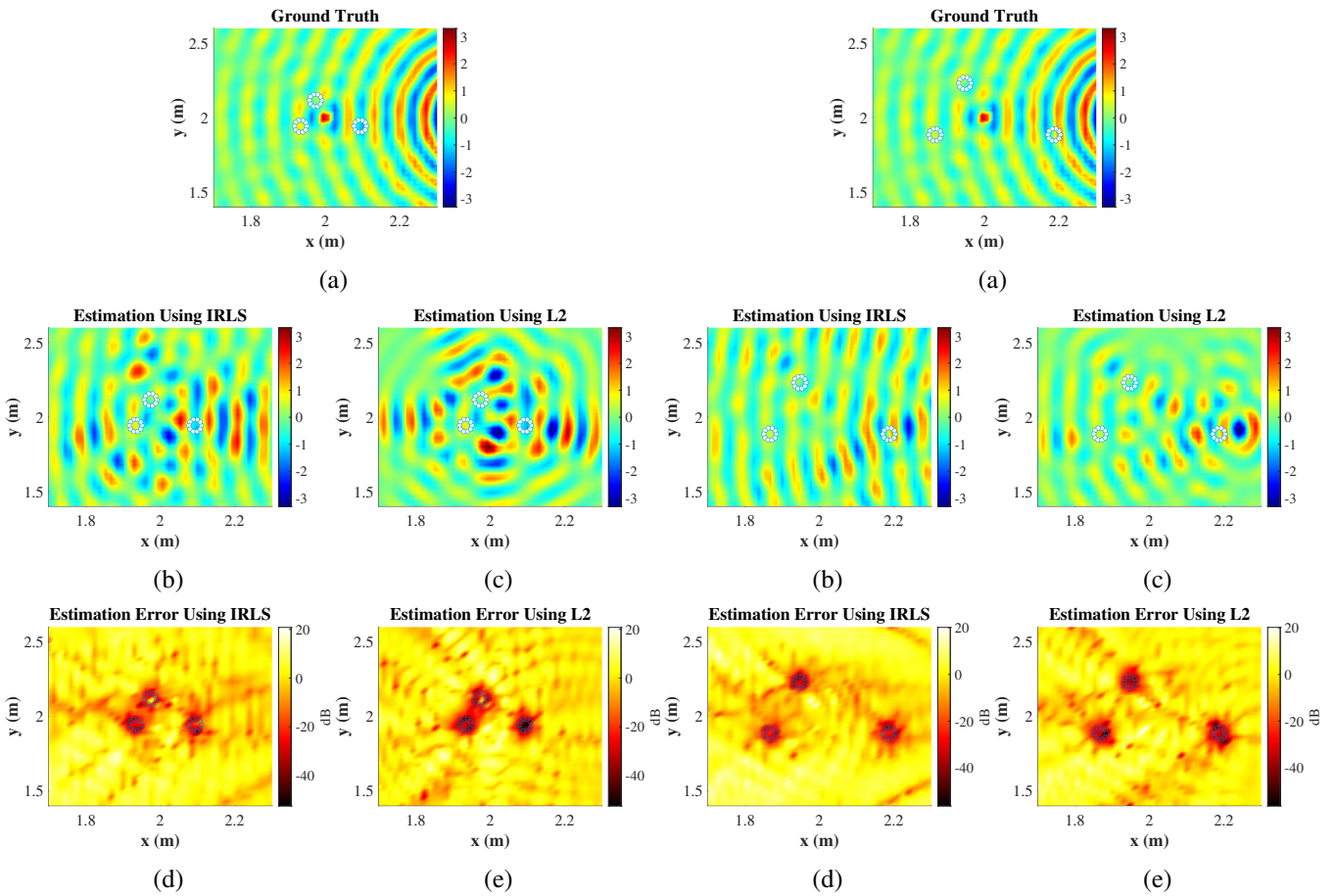


Figure 2: Simulation results for different estimation methods in case 1, where the distance between HOM1 and HOM2 is 0.4 m. (a) The ground truth of the original sound field. (b) Sound field estimation using IRLS sparse method. (c) Sound field estimation using ℓ_2 method. (d) Estimation error using IRLS sparse method. (e) Estimation error using ℓ_2 method.

Figure 3: Simulation results for different estimation methods in case 2, where the distance between HOM1 and HOM2 is 0.8 m. (a) The ground truth of the original sound field. (b) Sound field estimation using IRLS sparse method. (c) Sound field estimation using ℓ_2 method. (d) Estimation error using IRLS sparse method. (e) Estimation error using ℓ_2 method.

with an amplitude of 1 and another sound source locates outside the HOM array at [3.5, 2, 1.5] with an amplitude of 10. The two single-frequency signals are at 2000 Hz.

The virtual HOSs are placed on 3 planes, with two squared grids (50*50 HOSs) over the entire room above and below the x-y planes, respectively, and a circular grid (60 HOSs) outside the room on the x-y plane. The order of each virtual HOS is 2.

5.2 Results Analysis

We evaluate sound field estimation performance over ROI using ℓ_2 and IRLS methods under two different setups.

In case 1, from the comparison between Figure 2a and Figure 2b, we can see that the sound field estimation using IRLS sparse method can preserve the main directional pattern of the original sound field over the entire ROI. Whereas using ℓ_2 method can only preserve the original sound field at or near the three HOMs. Compared to Figure 2b, more artifacts can be observed at the boundary of the ROI in Figure 2c. This can be observed from the comparison between estimation errors over the ROI in Figure 2d and 2e as well.

We also investigate the effect of distance among HOMs to the accuracy of the soundfield estimation. Figure 3 has demonstrated the simulation results in case 2, where the distance among HOMs are relatively large. Compared with Figure 2, the estimation results in Figure 3 show more similarity to the ground truth. This indicates that in order to achieve accurate estimation over entire ROI, we need to carefully choose the distance among the recording devices. For large ROI and complex sound scenes, adding more HOMs to capture the sound scene is also a potential solution.

6. CONCLUSIONS

This paper proposed a virtual higher order source method for virtual navigation in the complex sound scenes. This method can combine recordings from multiple microphone arrays and decompose them into a grid of virtual higher order sources. The simulation results demonstrate the effectiveness in virtual navigation for a mixture of interior and exterior sound field. Further research directions include extensive perceptual studies and further applications in augmented reality.

7. REFERENCES

- [1] B. Rafaely, V. Tourbabin, E. Habets, Z. Ben-Hur, H. Lee, H. Gamper, L. Arbel, L. Birnie, T. Abhayapala, and P. Samarasinghe, "Spatial audio signal processing for binaural reproduction of recorded acoustic scenes—review and challenges," *Acta Acustica*, vol. 6, p. 47, 2022.
- [2] L. Birnie, T. Abhayapala, P. Samarasinghe, and V. Tourbabin, "Sound field translation methods for binaural reproduction," in *IEEE WASPAA*, pp. 140–144, IEEE, 2019.
- [3] M. Kentgens, A. Behler, and P. Jax, "Translation of a higher order ambisonic sound scene based on parametric decomposition," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 151–155, IEEE, 2020.
- [4] L. Birnie, T. Abhayapala, V. Tourbabin, and P. Samarasinghe, "Mixed source sound field translation for virtual binaural application with perceptual validation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 1188–1203, 2021.
- [5] L. Birnie, Z. Ben-Hur, V. Tourbabin, T. Abhayapala, and P. Samarasinghe, "Bilateral-ambisonic reproduction by soundfield translation," in *2022 International Workshop on Acoustic Signal Enhancement (IWAENC)*, pp. 1–5, IEEE, 2022.
- [6] E. Fernandez-Grande, D. Cavedes-Nozal, M. Hahmann, X. Karakonstantis, and S. A. Verbarg, "Reconstruction of room impulse responses over extended domains for navigable sound field reproduction," in *2021 Immersive and 3D Audio: from Architecture to Automotive (I3DA)*, pp. 1–8, IEEE, 2021.
- [7] L. McCormack, A. Politis, T. McKenzie, C. Hold, and V. Pulkki, "Object-based six-degrees-of-freedom rendering of sound scenes captured with multiple ambisonic receivers," *Journal of the Audio Engineering Society*, vol. 70, no. 5, pp. 355–372, 2022.
- [8] S. Damiano, F. Borra, A. Bernardini, F. Antonacci, and A. Sarti, "Soundfield reconstruction in reverberant rooms based on compressive sensing and image-source models of early reflections," in *2021 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pp. 366–370, IEEE, 2021.

- [9] X. Tang, J. Zhang, D. L. Alon, Z. Ben-Hur, P. Samarasinghe, and T. Abhayapala, "Wave domain sound field interpolation using two spherical microphone arrays," in *2022 30th European Signal Processing Conference (EUSIPCO)*, pp. 319–323, IEEE, 2022.
- [10] M. A. Poletti, T. Betlehem, and T. Abhayapala, "Higher order loudspeakers for improved surround sound reproduction in rooms," in *AES Convention 133*, 2012.
- [11] P. N. Samarasinghe, M. A. Poletti, S. A. Salehin, T. D. Abhayapala, and F. M. Fazi, "3d soundfield reproduction using higher order loudspeakers," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 306–310, IEEE, 2013.
- [12] J. A. Tropp and A. C. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Transactions on information theory*, vol. 53, no. 12, pp. 4655–4666, 2007.
- [13] S. Chen and D. Donoho, "Basis pursuit," in *Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers*, vol. 1, pp. 41–44, IEEE, 1994.
- [14] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [15] R. Chartrand and W. Yin, "Iteratively reweighted algorithms for compressive sensing," in *2008 IEEE international conference on acoustics, speech and signal processing*, pp. 3869–3872, IEEE, 2008.
- [16] "em32 eigenmike microphone array release notes," 2013.