

Sound Field Translation and Mixed Source Model for Virtual Applications with Perceptual Validation

Lachlan Birnie*[§], Thushara Abhayapala*, Vladimir Tourbabin[†], Prasanga Samarasinghe*

*Audio & Acoustic Signal Processing Group, The Australian National University, Canberra, Australia

[†]Facebook Reality Labs, Redmond, Washington, USA

Abstract—Non-interactive and linear experiences like cinema film offer high quality surround sound audio to enhance immersion, however the listener’s experience is usually fixed to a single acoustic perspective. With the rise of virtual reality, there is a demand for recording and recreating real-world experiences in a way that allows for the user to interact and move within the reproduction. Conventional sound field translation techniques take a recording and expand it into an equivalent environment of virtual sources. However, the finite sampling of a commercial higher order microphone produces an acoustic sweet-spot in the virtual reproduction. As a result, the technique remains to restrict the listener’s navigable region. In this paper, we propose a method for listener translation in an acoustic reproduction that incorporates a mixture of near-field and far-field sources in a sparsely expanded virtual environment. We perceptually validate the method through a Multiple Stimulus with Hidden Reference and Anchor (MUSHRA) experiment. Compared to the planewave benchmark, the proposed method offers both improved source localizability and robustness to spectral distortions at translated positions. A cross-examination with numerical simulations demonstrated that the sparse expansion relaxes the inherent sweet-spot constraint, leading to the improved localizability for sparse environments. Additionally, the proposed method is seen to better reproduce the intensity and binaural room impulse response spectra of near-field environments, further supporting the strong perceptual results.

Index Terms—Sound field navigation, translation, virtual-reality reproduction, binaural synthesis, MUSHRA, higher order microphone.

I. INTRODUCTION

Virtual reality devices will provide a novel framework for people to interact with each other at a higher social bandwidth through immersive audio and visual reproductions of the real-world [1], [2]. For example, in the future a person may be able to experience a live concert or orchestral performance through a virtual reproduction in their own home [3]. To complete the immersive experience, the listener/viewer should be allowed to explore and interact with the virtual reproduction [4]. Subsequently, methods to accurately record and model the perceptual change in visual and auditory information as the user translates are required to maintain the original experience.

Camera arrays have been used to capture visual information at multiple points-of-view for use in virtual reproductions [5]. Similarly, microphones distributed about an environment can record the spatial auditory scene from multiple points-of-view [6]. However, hardware and feasibility restrictions limit the

continuous space that can be recorded, and as a result, during reproduction the listener is usually stuck in the fixed acoustic perspective of the microphone [7].

Recently, there have been two key approaches towards extending the auditory range that a listener can navigate inside a virtual sound field reproduction. These are an interpolation-based [8] and an extrapolation-based approach [9]. Interpolation approaches utilize a grid of higher order microphones distributed about the acoustic space, and interpolate the sound to the listener during reproduction [10], [11]. Better coloration and localization performance is expected from interpolation than extrapolation [8]. However, interpolation may not be feasible for all real-world scenarios due to the large spatial, hardware, and synchronization costs associated with constructing a microphone grid [12]. Furthermore, typically listeners are confined to the interior region of the grid [13], and sound sources within the grid may cause comb-filtering spectral distortions [14]. Methods that alleviate these drawbacks and allow the listener to translate beyond the grid have been developed, however, they usually require additional localization and separation of direct sound field components [15], [16].

On the other hand, the extrapolation-based approach expands a single higher order microphone’s recording outwards to the translated listener’s position [17]. As a result, extrapolation overcomes many of the hardware and spatial drawbacks of the interpolation approach. Because a single microphone is utilized, the audio and visual capture system can occupy a single seat in the audience of a live event, which causes less obstruction and allows for more impromptu recordings.

Many extrapolation-based sound field translation methods have been developed to allow listener navigation in virtual reproduction, such as Ambisonic [17], [18], harmonic re-expansion [19], discrete source [20], and point-source distribution [21]. One of the most popular extrapolation-based methods which we consider to be the benchmark in this paper, is the planewave method [22]. In this method, the recorded acoustic environment is expanded into a secondary distribution of virtual planewave sources [23]. The secondary virtual environment constructs a sound field that is extrapolated from the microphone and is equivalent to the recording. The listener can perceptually move about the reproduction by translating the secondary environment’s sound field [20], [22].

In practice, however, most extrapolation-based approaches, including the planewave method, are constrained by the higher order microphone used for recording [17]. Hardware limitations result in an approximate and truncated sound field

[§]This research is supported by an Australian Government Research Training Program (RTP) Scholarship, and Facebook Reality Labs.

recording that is confined to a finite region [24], where the truncated recording is restricted by both the upper frequency band and the microphone radius [25]. As a result, the listener can only navigate within a small acoustic sweet-spot of a few centimeters which is defined by the commercial microphone's size [26]. Attempting to navigate beyond this inherent sweet-spot region, even after extrapolation, results in spectral distortions [27]–[29], degraded source localization [17], [30], and a poor perceptual listening experience.

Studies have shown that compensating for near-field effects attributes to better sound field reproduction [31]. Some reproduction methods have been able to model near-field sources with the use of prior knowledge of the source position or additional source localization processing [16], [32]–[34]. However, the planewave translation method is limited by its far-field virtual source model, which makes the reproduction of near-field propagation difficult [35].

In this paper we propose an alternative secondary source model for an extrapolated virtual reproduction that enables both a near-field and far-field propagation mixture. The method is built upon the benchmark planewave method which we review in Section II. We expand the truncated recording of a commercial higher order microphone [26], [36], [37] into a mixture of secondary virtual sources that are distributed in both the near-field and far-field (Section III), to create a more perceptually accurate reproduction. In addition to the source mixture, we also propose using a L1-norm regularization [38] to sparsely expand the recording into the equivalent virtual environment (Section III-C), as it has been shown to help extrapolate mode-limited sound fields [39]–[41].

We initially proposed the near-field far-field source mixture in [42] without any substantial verification of the method's perceptual performance. In this paper, we study the perceptual aspects of the source mixture through a perceptual listening test with human subjects and investigate the results against numerical simulations of the extrapolated pressure and intensity fields. The perceptual evaluation presented in Section IV, utilizes a MULTiple Stimulus with Hidden Reference and Anchor (MUSHRA) [43], [44] framework adapted for use in a virtual environment to provide listeners with both an auditory and visual reference of the real-world environment [12], [33]. We compare four translation methods with differing virtual source models and expansion techniques. We test the methods for the reproduction of human speech and music against the metrics of source localizability and robustness to spectral distortions. We will show that the proposed method offers greater perceptual accuracy and a more immersive experience for listeners moving throughout an expanded virtual reproduction. In Section V, we conduct a simulation study and show that the proposed method's perceptual performance is likely due to its ability to better reconstruct the near-field pressure and intensity of the original environment. We give our concluding remarks and suggestions for future work in Section VI.

II. PROBLEM FORMULATION AND THE PLANEWAVE SOUND FIELD TRANSLATION METHOD

In this section, we formulate the problem of reconstructing a recorded real-world experience such that a listener is able

to perceptually move through the acoustic reproduction. We first present the process of recording a general sound field with a commercial higher order microphone. We then review the planewave sound field translation method presented in [22], which segments the reproduction into three parts. First, a virtual acoustic environment is built from a superposition of planewave sources. Second, planewave driving signals are estimated from the recording to model an equivalent acoustic environment. Third, a listener is placed inside the virtual equivalent environment and binaural signals are rendered as they move. We provide a discussion on the perceptual shortcomings of this planewave translation method at the end.

A. Sound Field Capture

Consider a real-world acoustic environment that contains multiple sound sources, for example, a musical performance with many instruments. Let the origin \mathbf{o} denote the center of the environment's listening space, such as a seat in the middle of the audience. Each sound source is positioned at $\mathbf{z} = (r, \theta, \phi)$ with respect to \mathbf{o} , where $\theta \in [0, \pi]$ is the elevation angle downwards from the z-axis, and $\phi \in [0, 2\pi)$ is the azimuth angle counterclockwise from the x-axis. For a listener in the audience at position \mathbf{d} , the true sound they experience in the real-world can be described by

$${}^{\text{(real)}}P_{\{l,r\}}(k, \mathbf{d}) = \sum_{\mu=1}^U H_{\{l,r\}}(k, \mathbf{z}_\mu; \mathbf{d}) \times s_\mu(k), \quad (1)$$

where ${}^{\text{(real)}}P_{\{l,r\}}(k, \mathbf{d})$ is the pressure at the listener's left and right ear, $H_{\{l,r\}}(k, \mathbf{z}; \mathbf{d})$ is the transfer function between each source and the listener's ears, or simply the Head-Related Transfer Function (HRTF) when the listener is in a free-field space without any reflections, $s_\mu(k)$ is the sound signal of the μ^{th} source, $\mu = (1, \dots, U)$, $k = 2\pi f/c$ is the wave number, f is the frequency, and c is the speed of sound. From here on, we assume H to be the free-field HRTF for simplicity.

The aim is to record and reproduce the real-world auditory experience of (1) for every possible listening position. The homogeneous sound field that encompasses every arbitrary listening position \mathbf{x} , where $|\mathbf{x}| < |\mathbf{z}|$, can be expressed through a spherical harmonic decomposition of [45]

$$P(k, \mathbf{x}) = \sum_{n=0}^{\infty} \sum_{m=-n}^n \alpha_{nm}(k) j_n(k|\mathbf{x}|) Y_{nm}(\hat{\mathbf{x}}), \quad (2)$$

where $|\cdot| \equiv r$, $\hat{\cdot} \equiv (\theta, \phi)$, n and m are index terms denoting spherical harmonic order and mode, respectively, $j_n(\cdot)$ are the spherical Bessel functions of the first kind, $Y_{nm}(\cdot)$ are the set of spherical harmonic basis functions, and $\alpha_{nm}(k)$ are the sound field's coefficients which completely describe the source-free acoustic environment centered about \mathbf{o} when $\alpha_{nm}(k)$ is known for all $n \in [0, \infty)$.

In practice, the real-world acoustic environment can be recorded with an N^{th} order microphone, by estimating the sound field's $\alpha_{nm}(k)$ coefficients for a finite set of $n \in [0, N]$. Consider a N^{th} order microphone centered at \mathbf{o} , such as a spherical [26] (or planar [46], [47]) microphone array. The microphone array consists of $q = (1, \dots, Q)$ pressure sensors

that enclose the spherical acoustic region (listening space) of radius $|\mathbf{x}_Q|$ to be recorded. The sound field within this region can be estimated with [48]

$$\alpha_{nm}(k) \approx \sum_{q=1}^Q w_q \frac{P(k, \mathbf{x}_q) Y_{nm}^*(\hat{\mathbf{x}}_q)}{b_n(k|\mathbf{x}_Q|)}, \quad n \in [0, N], \quad (3)$$

where w_q are a set of suitable sampling weights [49], and $b_n(\cdot)$ is the rigid baffle equation [45].

However, commercial N^{th} order microphones can only record a small acoustic region ($|\mathbf{x}_Q| < 0.05$ m [26]) due to the hardware complexity and size constraint trade-offs with the spatial sampling Nyquist theorem [24]. The microphone's truncation order is restricted by the limited number of sensors, such that $Q \geq (N + 1)^2$. Furthermore, the microphone's recording region and frequency range are balanced by the $N = \lceil k|\mathbf{x}_Q| \rceil$ rule [50]. These two microphone properties define a maximum $|\mathbf{x}_Q|$ inside which the sound field is effectively of order $\leq N$. Beyond $|\mathbf{x}_Q|$, the reconstructed sound field requires higher orders $> N$ which are unknown, resulting in truncation error that degrades perceptual accuracy.

When reconstructing (1) from the recording, the left and right ear signals for the listener can be reassembled in the spherical harmonic domain by [51], [52]

$$P_{\{l,r\}}^{\text{(mic)}}(k, \mathbf{o}) = \sum_{n=0}^N \sum_{m=-n}^n H_{\{l,r\}}^{nm}(k) \times \alpha_{nm}(k), \quad (4)$$

where $H_{\{l,r\}}^{nm}(k)$ are the spherical harmonic decomposition coefficients of the HRTF $H_{\{l,r\}}(k, \mathbf{z}; \mathbf{o})$. In the reproduction (4), truncation forces the listener to the fixed auditory perspective of the microphone at \mathbf{o} . If the listener attempts to move then they would immediately translate beyond the $|\mathbf{x}_Q|$ boundary and begin to experience spectral distortions, degraded source localization performance, and a loss in perceptual immersion.

The objective of this paper is to relax this sweet-spot spatial constraint when reconstructing the sound field of a commercial microphone recording, and to build an equivalent virtual environment that allows for a listener to move about the acoustic space with a sustained perceptual immersion. For the remainder of this section we review the planewave sound field translation method that we consider to be the baseline method for enabling listener navigation.

B. Planewave Distribution

The planewave sound field translation method aims to construct a virtual acoustic environment that is perceptually equivalent to the real-world recording. The building block of this virtual environment is the planewave source, whose sound field is modeled as

$$P(k, \mathbf{x}) = \frac{e^{-ik\hat{\mathbf{y}} \cdot \mathbf{x}}}{4\pi}, \quad (5)$$

where $\hat{\mathbf{y}}$ denotes the planewave's incident direction. It is known that any acoustic free field can be modeled by an infinite superposition of planewaves [23]. Therefore, the equivalent virtual environment is constructed from a spherical distribution of virtual planewave sources, expressed as

$$P^{\text{(pw)}}(k, \mathbf{x}) = \int \psi(k, \hat{\mathbf{y}}; \mathbf{o}) \frac{e^{-ik\hat{\mathbf{y}} \cdot \mathbf{x}}}{4\pi} d\hat{\mathbf{y}}, \quad (6)$$

where $\psi(k, \hat{\mathbf{y}}; \mathbf{o})$ denotes the driving function of the planewave distribution as observed at \mathbf{o} . If the driving function is modeled correctly then the planewave distribution can recreate the acoustic environment, such that ${}^{\text{(pw)}}P(k, \mathbf{x}) = {}^{\text{(real)}}P(k, \mathbf{x})$. To achieve this, the driving function needs to be estimated/expanded from the recorded $\alpha_{nm}(k)$ coefficients, which we describe next.

C. Planewave Expansion

The sound field about \mathbf{o} due to a single virtual planewave can be expressed by the decomposition of [45]

$$\frac{e^{-ik\hat{\mathbf{y}} \cdot \mathbf{x}}}{4\pi} = \sum_{n=0}^{\infty} \sum_{m=-n}^n (-i)^n Y_{nm}^*(\hat{\mathbf{y}}) j_n(k|\mathbf{x}|) Y_{nm}(\hat{\mathbf{x}}). \quad (7)$$

Additionally, the driving function centered at \mathbf{o} can also be expressed in terms of a harmonic decomposition, given as

$$\psi(k, \hat{\mathbf{y}}; \mathbf{o}) = \sum_{n'=0}^{\infty} \sum_{m'=-n'}^{n'} \beta_{n'm'}(k) Y_{n'm'}(\hat{\mathbf{y}}), \quad (8)$$

where $\beta_{nm}(k)$ are the spherical harmonic decomposition coefficients of $\psi(k, \hat{\mathbf{y}}; \mathbf{o})$ which describe the sound field about the planewave distribution. Substituting both (7) and (8) into (6) gives the planewave distribution's sound field in spherical harmonics, as

$$P^{\text{(pw)}}(k, \mathbf{x}) = \sum_{n=0}^{\infty} \sum_{m=-n}^n \underbrace{(-i)^n \beta_{nm}(k)}_{\alpha_{nm}(k)} j_n(k|\mathbf{x}|) Y_{nm}(\hat{\mathbf{x}}). \quad (9)$$

From (9), the relationship between the $\beta_{nm}(k)$ coefficients and the recorded $\alpha_{nm}(k)$ coefficients can be extracted. Rearranging this relationship for $\beta_{nm}(k) = i^n \alpha_{nm}(k)$, expresses a planewave distribution that is equivalent to the recorded environment. Substituting this relationship back into (8), gives a closed-form expansion for a planewave driving function that matches the recording,

$$\psi(k, \hat{\mathbf{y}}; \mathbf{o}) = \sum_{n=0}^N \sum_{m=-n}^n i^n \alpha_{nm}(k) Y_{nm}(\hat{\mathbf{y}}). \quad (10)$$

Synthesizing a virtual environment with this driving function through (6) produces a sound field that is equivalent to the recording. However, the recording (3) is only an approximation of the real environment, and therefore (10) is also an approximate, such that ${}^{\text{(pw)}}P(k, \mathbf{x}) \equiv {}^{\text{(mic)}}P(k, \mathbf{x}) \approx {}^{\text{(real)}}P(k, \mathbf{x})$.

D. Planewave Auralization

A listener inside the planewave distribution is immersed within a spatial reproduction of the real-world acoustic environment. The binaural signals for the listener at the distribution center can be presented by exchanging their HRTF into (6), giving

$$P_{\{l,r\}}^{\text{(pw)}}(k, \mathbf{o}) = \int \psi(k, \hat{\mathbf{y}}; \mathbf{o}) H_{\{l,r\}}(k, \hat{\mathbf{y}}; \mathbf{o}) d\hat{\mathbf{y}}. \quad (11)$$

Furthermore, the planewave distribution allows for the listener to move perceptually about the reproduction. The sound heard

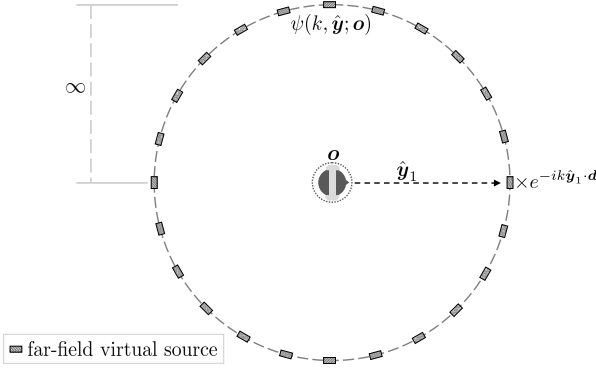


Fig. 1: Illustration of the equivalent virtual planewave distribution. The listener's perspective is fixed at the distribution center \mathbf{o} , where a phase shift applied to the driving function translates the sound field about the listener.

by the listener who is translated to $\mathbf{x} = [\mathbf{o} + \mathbf{d}] \equiv \mathbf{d}$ can be derived from (6) as

$$\begin{aligned} {}^{(\text{pw})}P(k, [\mathbf{o} + \mathbf{d}]) &= \int \psi(k, \hat{\mathbf{y}}; \mathbf{o}) \frac{e^{-ik\hat{\mathbf{y}} \cdot [\mathbf{o} + \mathbf{d}]}}{4\pi} d\hat{\mathbf{y}} \\ &= \int \psi(k, \hat{\mathbf{y}}; \mathbf{o}) e^{-ik\hat{\mathbf{y}} \cdot \mathbf{d}} \frac{e^{-ik\hat{\mathbf{y}} \cdot \mathbf{o}}}{4\pi} d\hat{\mathbf{y}}. \end{aligned} \quad (12)$$

It is observed from (12) that the translation in space differs only by a phase shift in the planewave driving function. Therefore, applying the translational phase shift of [35]

$$\psi(k, \hat{\mathbf{y}}; \mathbf{d}) = \psi(k, \hat{\mathbf{y}}; \mathbf{o}) \times e^{-ik\hat{\mathbf{y}} \cdot \mathbf{d}}, \quad (13)$$

to the binaural signals in (11), allows for the listener to dynamically move their acoustic perspective by

$${}^{(\text{pw})}_{\{l,r\}}P(k, \mathbf{d}) = \int \psi(k, \hat{\mathbf{y}}; \mathbf{d}) H_{\{l,r\}}(k, \hat{\mathbf{y}}; \mathbf{o}) d\hat{\mathbf{y}}. \quad (14)$$

In practice, the virtual planewave distribution (6) can be realized with a discrete set of known sources,

$${}^{(\text{pw})}P(k, \mathbf{x}) \approx \sum_{\ell=1}^L w_{\ell} \psi(k, \hat{\mathbf{y}}_{\ell}; \mathbf{o}) \frac{e^{-ik\hat{\mathbf{y}}_{\ell} \cdot \mathbf{x}}}{4\pi} \quad (15)$$

where $\ell = (1, \dots, L)$ index each virtual planewave, L is the total number of sources, and w_{ℓ} are a set of suitable sampling weights. Similarly, the dynamic binaural signals can be realized from the discrete distribution with

$${}^{(\text{pw})}_{\{l,r\}}P(k, \mathbf{d}) \approx \sum_{\ell=1}^L w_{\ell} \psi(k, \hat{\mathbf{y}}_{\ell}; \mathbf{d}) H_{\{l,r\}}(k, \hat{\mathbf{y}}_{\ell}; \mathbf{o}). \quad (16)$$

We illustrate this planewave method to sound field translation in Fig. 1. The reproduction is expressed by many discrete planewave signals that are known continuously throughout the virtual environment. Therefore, the planewave method does not explicitly limit the amount the listener can translate. However, (16) uses $\psi(k, \hat{\mathbf{y}}; \mathbf{o})$ which is estimated through (3) and (10). As a result, the N^{th} order truncation inherently remains, and the listener's movement is still implicitly limited.

E. Discussion

The virtual planewave expansion enables listener translation, however, some shortcomings are still exhibited in the listener's perception:

- As mentioned, the planewave method inherits truncation artifacts through an over-approximation of (10), and the listener's movement remains inherently restricted inside the virtual reproduction. As the listener translates further away from the recording's sweet-spot, they begin to experience spectral distortions, a loss in source localization, and poorer perceptual accuracy.
- The planewave expansion has difficulties in synthesizing near-field sound sources due to its far-field source model.
- The planewave auralization (16) fixes the HRTF perspective to the virtual distribution center, $H_{\{l,r\}}(k, \hat{\mathbf{y}}_{\ell}; \mathbf{o})$, and performs translation by phase shifting the sound field with (13). However, the HRTF propagation vectors $\hat{\mathbf{y}}_{\ell}; \mathbf{o}$ remain un-shifted, and as a result, the HRTF models untranslated head reflections as the listener moves.

In the next section we propose an alternative sound field translation model to address the above shortcomings.

III. MIXEDWAVE SOUND FIELD TRANSLATION METHOD

In this section, we define a virtual source that models both a near-field and far-field propagation, which we will refer to as a mixedwave source. We then build a virtual distribution of mixedwave sources and expand a real-world recording into an equivalent sound field. Additionally, we also propose a sparse method for expanding a virtual source distribution that alleviates some of the spatial restrictions imposed by the truncated recording.

A. Mixture of Near-Field and Far-Field Sources

Here, we define the virtual source that will be the building block for our proposed method. Consider a near-field point-source at \mathbf{y} , where the driving signal of the source with respect to itself is denoted $\dot{\psi}(k, \mathbf{y})$. We can express the driving function observed at a position \mathbf{x} with [45]

$$\psi(k, \mathbf{y}; \mathbf{x}) = \dot{\psi}(k, \mathbf{y}) \frac{e^{ik\|\mathbf{y}-\mathbf{x}\|}}{\|\mathbf{y}-\mathbf{x}\|}. \quad (17)$$

Evaluating (17) when $\mathbf{x} = \mathbf{o}$ gives the driving function observed by a receiver/microphone, as

$$\psi(k, \mathbf{y}; \mathbf{o}) = \dot{\psi}(k, \mathbf{y}) \frac{e^{ik\|\mathbf{y}\|}}{\|\mathbf{y}\|}. \quad (18)$$

Rearranging (18) gives an expression for the source signal in terms of the source's distance and the driving function observed by the receiver,

$$\dot{\psi}(k, \mathbf{y}) = \psi(k, \mathbf{y}; \mathbf{o}) \|\mathbf{y}\| e^{-ik\|\mathbf{y}\|}. \quad (19)$$

Substituting (19) back into (17) provides the driving function observed at any arbitrary point \mathbf{x} in terms of the function observed by the receiver/microphone, expressed as

$$\psi(k, \mathbf{y}; \mathbf{x}) = \underbrace{\psi(k, \mathbf{y}; \mathbf{o}) \|\mathbf{y}\|}_{\dot{\psi}(k, \mathbf{y})} e^{-ik\|\mathbf{y}\|} \frac{e^{ik\|\mathbf{y}-\mathbf{x}\|}}{\|\mathbf{y}-\mathbf{x}\|}. \quad (20)$$

We note that the $|\mathbf{y}|e^{-ik|\mathbf{y}|}$ term can be seen to have redefined the point-source from being a function with respect to itself, to being a function with respect to \mathbf{o} . This allows us to observe the source distribution at \mathbf{o} with a microphone and estimate the sound at any translated position \mathbf{x} .

Additionally, the constant term has the property of [25]

$$\lim_{|\mathbf{y}| \rightarrow \infty} |\mathbf{y}|e^{-ik|\mathbf{y}|} \frac{e^{ik|\mathbf{y}-\mathbf{x}|}}{4\pi|\mathbf{y}-\mathbf{x}|} = \frac{e^{-ik\hat{\mathbf{y}} \cdot \mathbf{x}}}{4\pi}, \quad (21)$$

which allows for a mixture of near-field and far-field virtual source distributions to be modeled with this building block. We define this building block as the mixedwave source,

$$P(k, \mathbf{x}) = |\mathbf{y}|e^{-ik|\mathbf{y}|} \frac{e^{ik|\mathbf{y}-\mathbf{x}|}}{4\pi|\mathbf{y}-\mathbf{x}|}. \quad (22)$$

In the spherical harmonic domain,

$$|\mathbf{y}|e^{-ik|\mathbf{y}|} \frac{e^{ik|\mathbf{y}-\mathbf{x}|}}{4\pi|\mathbf{y}-\mathbf{x}|} = \sum_{n=0}^{\infty} \sum_{m=-n}^n ik|\mathbf{y}|e^{-ik|\mathbf{y}|} h_n(k|\mathbf{y}|) Y_{nm}^*(\hat{\mathbf{y}}) j_n(k|\mathbf{x}|) Y_{nm}(\hat{\mathbf{x}}), \quad (23)$$

where $h_n(\cdot)$ is the spherical Hankel function of the first kind. We note that spherical Hankel functions also have

$$\lim_{|\mathbf{y}| \rightarrow \infty} ik|\mathbf{y}|e^{-ik|\mathbf{y}|} h_n(k|\mathbf{y}|) = (-i)^n, \quad (24)$$

to correspond with (21). We can observe from (24) that when the mixedwave source is placed in the far-field, the definition of (23) will match that of the planewave source (7). This property then allows for both a near-field sound propagation to be modeled by a mixedwave distribution with a small radius, and a far-field sound propagation modeled by a mixedwave distribution with a large radius. We will use this near-field and far-field distribution of mixedwave sources as the basis of our proposed sound field translation method next.

B. Mixedwave Method for Sound Field Translation

Following the planewave translation method, our proposed mixedwave translation method is also broken into three parts.

1) *Mixedwave Distribution*: We propose constructing a virtual equivalent sound field from two concentric spherical distributions of mixedwave sources. The first virtual sphere is placed in the near-field with a radius of $R_{(\text{nf})}$, and the second sphere is placed at $R_{(\text{ff})}$ in the far-field, such that

$$\begin{aligned} &^{(\text{mw})}P(k, \mathbf{x}) \\ &= \int \psi(k, R_{(\text{nf})}\hat{\mathbf{y}}; \mathbf{o}) R_{(\text{nf})} e^{-ikR_{(\text{nf})}} \frac{e^{ik|R_{(\text{nf})}\hat{\mathbf{y}}-\mathbf{x}|}}{4\pi|R_{(\text{nf})}\hat{\mathbf{y}}-\mathbf{x}|} d\hat{\mathbf{y}} \\ &+ \int \psi(k, R_{(\text{ff})}\hat{\mathbf{y}}; \mathbf{o}) R_{(\text{ff})} e^{-ikR_{(\text{ff})}} \frac{e^{ik|R_{(\text{ff})}\hat{\mathbf{y}}-\mathbf{x}|}}{4\pi|R_{(\text{ff})}\hat{\mathbf{y}}-\mathbf{x}|} d\hat{\mathbf{y}}, \end{aligned} \quad (25)$$

where, $\psi(k, R\hat{\mathbf{y}}; \mathbf{o})$, $R \in \{R_{(\text{nf})}, R_{(\text{ff})}\}$, are the driving functions of the two mixedwave distributions centered at \mathbf{o} .

2) *Mixedwave Expansion*: Following the procedure in Section II-C, we can decompose the $\psi(k, R\hat{\mathbf{y}}; \mathbf{o})$ driving function into spherical harmonic aperture coefficients of $\beta_{n'm'}(k, R)$, expressed as

$$\psi(k, R\hat{\mathbf{y}}; \mathbf{o}) = \sum_{n'=0}^{\infty} \sum_{m'=-n'}^{n'} \beta_{n'm'}(k, R) Y_{n'm'}(\hat{\mathbf{y}}). \quad (26)$$

We substitute both (26) and (23) into (25) to extract the relationship between $\beta_{nm}(k)$ and $\alpha_{nm}(k)$, given as

$$\beta_{nm}(k, R) = \sum_{n=0}^N \sum_{m=-n}^n \frac{\alpha_{nm}(k)}{ikR e^{-ikR} h_n(kR)}. \quad (27)$$

Finally, we substitute (27) back into (26) to derive a closed-form expansion for the mixedwave driving functions in terms of the recorded coefficients,

$$\psi(k, R\hat{\mathbf{y}}; \mathbf{o}) = \sum_{n=0}^N \sum_{m=-n}^n \frac{\alpha_{nm}(k)}{ikR e^{-ikR} h_n(kR)} Y_{nm}(\hat{\mathbf{y}}). \quad (28)$$

We use a set of real-world recorded coefficients $\alpha_{nm}(k)$ with (28) to estimate the driving functions of the near-field and far-field virtual distributions, such that $^{(\text{mw})}P(k, \mathbf{x}) \equiv ^{(\text{mic})}P(k, \mathbf{x}) \approx ^{(\text{real})}P(k, \mathbf{x})$.

3) *Mixedwave Auralization*: Consider a listener inside the virtual mixedwave distribution at the translated position $\mathbf{x} = [\mathbf{o} + \mathbf{d}] \equiv \mathbf{d}$, $|\mathbf{d}| < R_{(\text{nf})}$, as shown in Fig. 2. We render the left and right binaural signals by applying the mixedwave driving function to the HRTF based on the listener's translated position, given as [17]

$$^{(\text{mw})}_{\{\text{l,r}\}}P(k, \mathbf{d}) = \int \psi(k, R\hat{\mathbf{y}}; \mathbf{o}) H_{\{\text{l,r}\}}(k, R\hat{\mathbf{y}}; \mathbf{d}) d\mathbf{y}, \quad (29)$$

where $R\hat{\mathbf{y}}; \mathbf{d}$ denotes the propagation direction of the mixedwave source with respect to \mathbf{d} , which is given by $(\mathbf{y} - \mathbf{d})$. We note that this is possible for the mixedwave distribution due to the finite positions of each source, unlike the infinite definitions of planewave sources.

Once again, a set of discrete sources can be used to practically realize the virtual mixedwave distributions, expressed as $^{(\text{mw})}P(k, \mathbf{x}) \approx$

$$\begin{aligned} &\sum_{\ell=1}^L w_{\ell} \psi(k, R_{(\text{nf})}\hat{\mathbf{y}}_{\ell}; \mathbf{o}) R_{(\text{nf})} e^{-ikR_{(\text{nf})}} \frac{e^{ik|R_{(\text{nf})}\hat{\mathbf{y}}_{\ell}-\mathbf{x}|}}{4\pi|R_{(\text{nf})}\hat{\mathbf{y}}_{\ell}-\mathbf{x}|} \\ &+ \sum_{\ell=1}^L w_{\ell} \psi(k, R_{(\text{ff})}\hat{\mathbf{y}}_{\ell}; \mathbf{o}) R_{(\text{ff})} e^{-ikR_{(\text{ff})}} \frac{e^{ik|R_{(\text{ff})}\hat{\mathbf{y}}_{\ell}-\mathbf{x}|}}{4\pi|R_{(\text{ff})}\hat{\mathbf{y}}_{\ell}-\mathbf{x}|}, \end{aligned} \quad (30)$$

where the near-field and far-field distributions each contain L sources. Similarly, we realize the mixedwave auralization within the discrete virtual distributions by

$$^{(\text{mw})}_{\{\text{l,r}\}}P(k, \mathbf{d}) = \sum_{\ell=1}^{2L} w_{\ell} \psi(k, \mathbf{y}_{\ell}; \mathbf{o}) H_{\{\text{l,r}\}}(k, \mathbf{y}_{\ell}; \mathbf{d}), \quad (31)$$

where $|\mathbf{y}_{\ell}| = R_{(\text{nf})}$ for $\ell \in [1, L]$, and $|\mathbf{y}_{\ell}| = R_{(\text{ff})}$ for $\ell \in [L+1, 2L]$, and $\mathbf{y}_{\ell}; \mathbf{d}$ is the propagation direction of the ℓ^{th} mixedwave source with respect to the translated listener. Unlike the planewave method, the maximum distance a listener can translate within the mixedwave environment is restricted

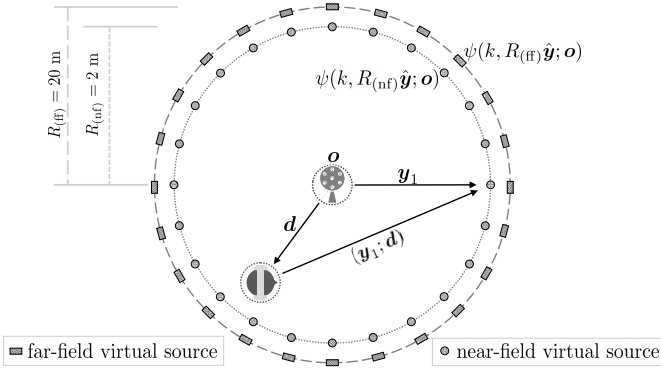


Fig. 2: Illustration of the equivalent virtual mixedwave sound field. The listener is translated to \mathbf{d} , and the vectors $(\mathbf{y}_\ell; \mathbf{d})$ are updated with the HRTF to auralize an immersive reproduction.

by $R_{(\text{nf})}$. However, we suspect that $R_{(\text{nf})}$ can be selected to match the size of a small real-world room that is recorded.

C. Sparse Expansion Methods

The closed-form expansion constructs a virtual environment that is equivalent to the original recording. However, the expansion distributes energy $\psi(k, \mathbf{y}_\ell; \mathbf{o})$ throughout all virtual sources. This causes an over-approximation of the truncated recording's underlying spatial artifacts. As a result, the amount a listener can translate before experiencing a loss in immersion is still inherently restricted by the recording's truncation. Furthermore, it is believed that modeling fewer virtual sources from propagation directions that are similar to the original environment will lead to better perceptual immersion [42]. For these reasons, we propose a sparse constrained expansion method for constructing our virtual mixedwave environment.

The coefficients $\alpha_{nm}(k)$ observed at the center of a virtual distribution can be expressed in matrix form as

$$\mathbf{A}\boldsymbol{\psi} = \boldsymbol{\alpha}, \quad (32)$$

where $\boldsymbol{\alpha} = [\alpha_{00}(k), \alpha_{1-1}(k), \dots, \alpha_{NN}(k)]^T$ are the recorded coefficients, $\boldsymbol{\psi} = [\psi(k, \mathbf{y}_1; \mathbf{o}), \dots, \psi(k, \mathbf{y}_L; \mathbf{o})]$ are the \mathcal{L} equivalent virtual source driving signals, and \mathbf{A} is the $(N+1)^2$ by \mathcal{L} expansion matrix. The entries of \mathbf{A} are given by $(-i)^n Y_{nm}^*(\hat{\mathbf{y}}_\ell)$ for a planewave expansion (10), and $ik|\mathbf{y}_\ell|e^{-ik|\mathbf{y}_\ell|}h_n(k|\mathbf{y}_\ell)Y_{nm}^*(\hat{\mathbf{y}}_\ell)$ where $\mathcal{L} = 2L$ for the two source distributions of a mixedwave expansion (28). We assume $L > (N+1)^2$ for the under-determined case.

We construct a sparse source distribution by solving the linear regression problem (32) using Iteratively Reweighted Least Squares (IRLS) [38]. In brief, the IRLS approach replaces the ℓ^p -objective function

$$\min_{\boldsymbol{\psi}} \|\boldsymbol{\psi}\|_p^p \quad \text{subject to } \mathbf{A}\boldsymbol{\psi} = \boldsymbol{\alpha}, \quad (33)$$

with a weighted ℓ^2 -norm,

$$\min_{\boldsymbol{\psi}} \sum_{i=1}^{\mathcal{L}} w_i \psi_i^2 \quad \text{subject to } \mathbf{A}\boldsymbol{\psi} = \boldsymbol{\alpha}, \quad (34)$$

where $w_i = |\psi_i^{(\nu-1)}|^{p-2}$ are the weights computed from the previous iterate $\boldsymbol{\psi}^{(\nu-1)}$. The next iterate is given by

$$\boldsymbol{\psi}^{(\nu)} = \mathbf{Q}_\nu \mathbf{A}^T (\mathbf{A} \mathbf{Q}_\nu \mathbf{A}^T)^{-1} \boldsymbol{\alpha}, \quad (35)$$

where \mathbf{Q}_ν is the diagonal matrix with $1/w_i = |\psi_i^{(\nu-1)}|^{2-p}$. Other regularization techniques can also be utilized, such as the Least-Absolute Shrinkage and Selection Operator (Lasso) [53], [54], and we direct the reader to [55] for further information in regards to compressive sensing.

D. Discussion

Continuing our discussion on the planewave method's shortcomings in Sec. II-E, we give the following comments:

- Sparsely expanding the virtual source distribution (32) with IRLS is expected to further enhance the perceptual immersion for a listener, as they should experience more localized virtual sources. Additionally, the sparsity relaxes the spatial sweet-spot restriction and over-approximation issue stemming from the closed-form expansion used by the planewave method. These properties are demonstrated by experiment in Section IV and by simulation in Section V.
- The mixedwave distribution can easily synthesize near-field sound sources. The modified point-source (22) can model a spherical-wave propagation by simply positioning the mixedwave source in the near-field.
- The mixedwave auralization (31) translates the HRTF with the listener. As a result, the propagation vectors in $H_{\{\text{l},\text{r}\}}(k, \mathbf{y}_\ell; \mathbf{d})$ are updated with \mathbf{d} to render changes in head reflection, similar to virtual higher-order Ambisonics [17]. Intuitively, this is expected to result in greater perceptual immersion.

We examine the perceptual advantages of the sparse expansion and the mixedwave source model against the planewave benchmark experimentally in the next section.

IV. PERCEPTUAL EXPERIMENT

Our aim is to maintain the immersion for a listener inside an acoustic reproduction. Therefore, it is of crucial importance, foremost, that we evaluate the proposed method against the planewave benchmark in a perceptual listening experiment. This section outlines the perceptual experiment system we implemented and presents the statistical results at the end.

A. Experiment Methodology

1) *Compared Methods*: We conducted a MUSHRA perceptual experiment to compare four translation methods. In total the experiment presented six signals:

- *Reference / hidden reference*: Signals of the true free-field transfer function between a real-world point-source and the translated listener, given by (1).
- *Anchor*: Signals of the truncated recording that is fixed spatially to the microphone's position (4). Sound field rotation is still rendered, but no translation is processed. This is a similar anchor to the three-degrees-of-freedom used in [33].

- *Benchmark / planewave closed-form (PW-CF)*: Signals rendered of a virtual planewave distribution (16) that are expanded through the closed-form expression (10).
- *Planewave IRLS (PW-IRLS)*: Signals of a IRLS (Section III-C) sparsely expanded planewave distribution.
- *Mixedwave closed-form (MW-CF)*: Signals rendered of a virtual mixedwave distribution (31) that are expanded through the closed-form expression (28).
- *Proposed method / mixedwave IRLS (MW-IRLS)*: Signals of a IRLS sparsely expanded mixedwave distribution.

The experiment comprised of four tests with two scoring metrics, *source localization* and *basic audio quality*, and two sound-signals, *speech* and *music*. The source localization test asked listeners to score on the perceived direction of the sound-source, the source width, and the sound field sparseness with respect to both a visual-reference and the reference signal. Whereas, the basic audio quality test asked listeners to score against the reference for spectral distortions and other audible processing artifacts. In total the scores of 17 participants were collected for the speech sound-source, and 11 scores for the music sound-source. The recording microphone was shown in the virtual environment, and listeners were informed that the further they translate, the greater the differences they should perceive between methods. We asked the listeners to score while accounting for each method’s performance over a 1 m square reproduction space.

2) *Experiment System*: We used an Oculus Rift along with a pair of Beyerdynamic DT 770 pro headphones to track the listener and provide a visual reference of the true sound source. We used the HRTFs of the FABIAN head and torso simulator [56] from the HUTUBS dataset [57], [58] for auralization. The HRTFs were rotated for each test signal by multiplying the HRTF coefficients with Wigner-D functions [59]. Signals were processed at a frame size of 4096 with 50% overlap and a 16 kHz sampling frequency, due to hardware constraints and the computational costs of the real-time experiment.

3) *Virtual Environment*: We simulated the real-world auditory experience with a single free-field point-source in order to generate a true experiment reference signal for the listener at every position. We constructed a virtual environment with \mathbf{o} placed at the center, and the XY-plane 1.25 m above the ground to align with a listener’s head while sitting. We modeled the *true* sound-source with a static point-source at (1, 0, 0) m. By *true*, we signify that the sound field generated by this point-source is denoted as the real-world auditory experience we record and reproduce.

Additionally, we also simulated the process of recording the truncated sound field of the true point-source. We used a 4th order rigid 36-sensor spherical microphone array centered at \mathbf{o} . Microphone sensors were distributed at Fliege positions [60] with 0.042 m radius to best represent a commercial microphone [26]. Recordings were generated by convolving the sound-source’s signal with the microphone’s impulse response. The $\alpha_{nm}(k)$ coefficients were extracted with (3) before being expanded into virtual distributions.

The planewave distribution consisted of $L = 36$ virtual sources at Fliege positions [60]. This selection was made as a trade-off with computation complexity. However, adding more

planewaves is not expected to improve source localization performance, as the distribution already over-samples the 4th order recording [9], [28], [30]. Similarly, the mixedwave distribution consisted of two sets of $L = 36$ virtual sources at the same Fliege positions. The first set was distributed in the near-field at $R_{(\text{nf})} = 2$ m, and the second was placed at $R_{(\text{ff})} = 20$ m in the far-field.

4) *Experiment Auralization*: The reference was rendered by convolving (in frequency domain) the sound-source signal with the true source-to-listener HRTF. For the anchor (4), the signals were convolved by multiplying $\alpha_{nm}(k)$ with the spherical HRTF-coefficients directly [52]. The planewave method signals were rendered with the convolution of the HRTF at \mathbf{o} and the phase-shifted driving function (16). The phase-shift was updated with the Oculus head position to render perceptual translation. For the mixedwave methods, the HRTFs were reconstructed between each source and the listener’s translated position. Binaural signals were then rendered with the convolution of the mixedwave driving function and the \mathbf{y}_ℓ -to- \mathbf{d} HRTF (31).

B. Experiment Results

1) *Box Plot*: Figure 3 shows the perceptual scores of the translation methods across all four tests. A Lilliefors test ($p_{\text{val}} > 0.01$) showed that our collected scores met the requirement for normal distribution, and a Tukey-Kramer multiple comparison test with 95% confidence was used to determine statistical significance. We discuss the results of these scores through an analysis of variance (ANOVA) examination.

2) *One-factor ANOVA results*: We used a one-factor ANOVA to determine if any of the translation methods performed significantly different in each of the perception tests. For speech localization (Fig. 3a), both MW-CF and MW-IRLS showed a significant improvement in score ($F_{(3,64)} = 6.2, p < 0.001$) compared to the PW-CF benchmark. Similar results ($F_{(3,64)} = 16.25, p < 0.001$) are shown for speech quality (Fig. 3b), where the mixedwave methods were found to be significantly different to PW-IRLS in addition to the benchmark. In the music sound-source tests (Fig. 3c and Fig. 3d), only MW-IRLS showed significantly improved means over the benchmark, while MW-CF did not. However, MW-CF was still observed to perform well for music localization in Fig. 3c and music quality in Fig. 3d as indicated by the significant median scores.

3) *Two-factor ANOVA results*: We performed a two-factor ANOVA to compare the effects of source-type (planewave and mixedwave) and expansion-type (closed-form and IRLS). In all four tests ($p \leq 0.008$), mixedwave source distributions were found to score higher means than planewave distributions. Whereas, a significant difference in expansion-type was only found in the speech sound-source tests, with IRLS showing better scores. For music localization ($F_{(1,40)} = 3.36, p = 0.074$) and music quality ($F_{(1,40)} = 1.07, p = 0.307$), no significant difference was found between closed-form and sparse expansions. Lastly, no interaction effects ($p \geq 0.382$) between virtual source-type and expansion-type were found.

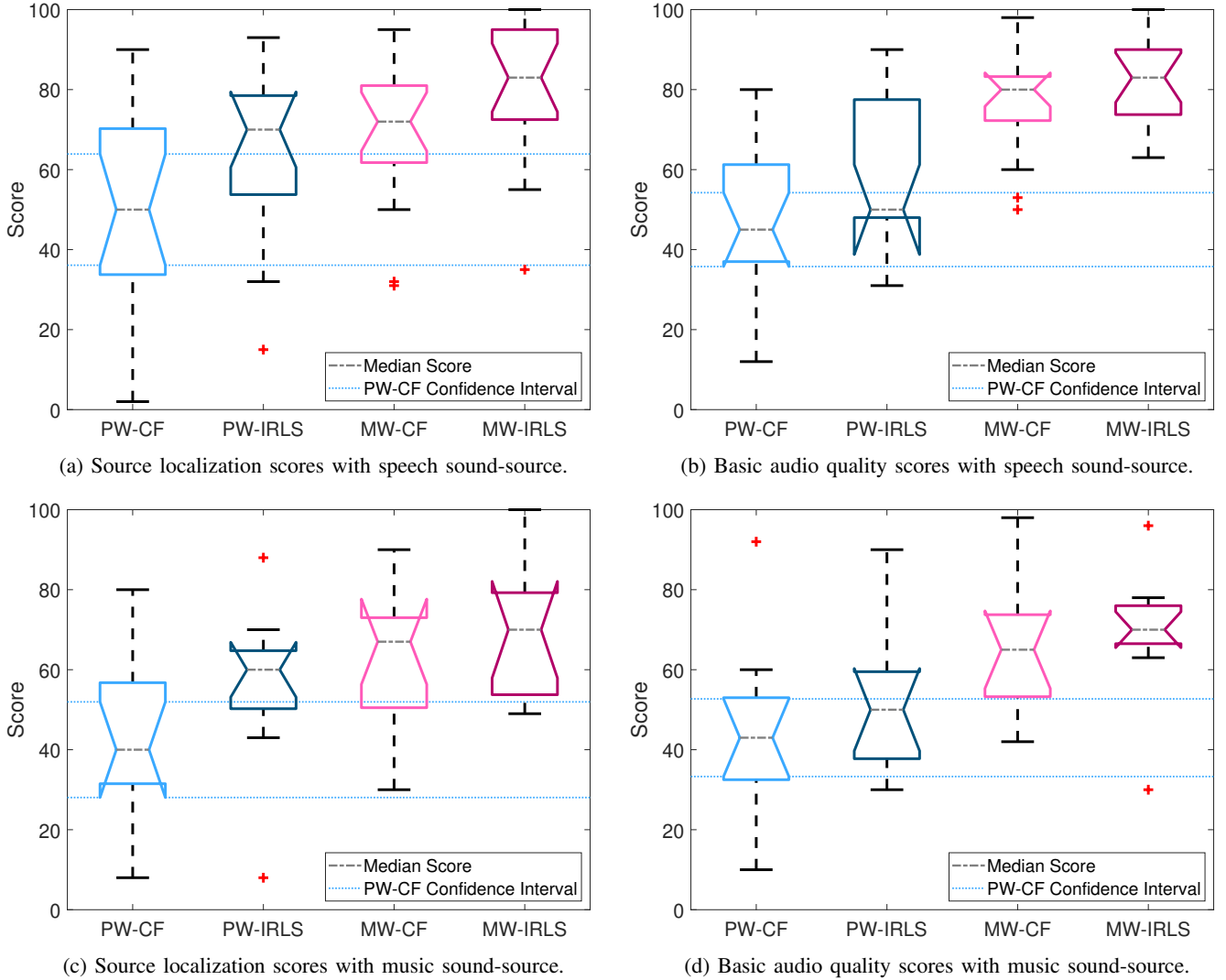


Fig. 3: Box plot of perception experiment scores for source localization (a) and (c), and basic audio quality (b) and (d). Each box bounds the interquartile range (IQR) with the center bar indicating the median score, and the whiskers extended to a maximum of $1.5 \times \text{IQR}$. The v -shaped notches in the box refer to the 95% confidence interval. When the notches between two boxes do not overlap, it can be concluded with 95% confidence that the true medians differ.

4) *Summary and discussion:* The proposed MW-IRLS method showed an improvement against the PW-CF benchmark in the perceptual criteria of source localization and audio quality for both a speech and music source. Furthermore, MW-CF also received higher mean scores when reconstructing human speech, and higher median scores for music. When comparing virtual expansion-types, the IRLS expansion was seen to have better quality robustness and localizability for a speech source, but not a music source. This may be explained by the IRLS matching the sparseness of the single human’s speech, but not the natural sound of music which is normally generated by multiple sound-sources. Nonetheless, this paper focuses on the modeling of secondary virtual sources. No interaction effect between the source model and expansion-type was found. This indicates that the strong perceptual results achieved by mixedwave methods were not dependent on the expansion-type, and are instead an outcome of the near-

field and far-field virtual source mixture. In the next section, we conduct a simulation analysis on the sound fields used in this experiment to gain further insight on properties that may have influenced these strong perceptual results.

V. SIMULATION ANALYSIS

In this section we simulate the same virtual environments that were used in the perception test (Section IV-A3). We examine their pressure and intensity fields to identify factors that may correlate with perceptual performance.

A. Error Metrics

We define the pressure error (PE) and intensity magnitude error (IME) between the true and reproduced sound field as

$$\left(\text{PE} = \frac{|P - \tilde{P}|^2}{|P|^2}, \quad \text{IME} = \frac{\|I - \tilde{I}\|^2}{\|I\|^2} \right) \times 100(\%). \quad (36)$$

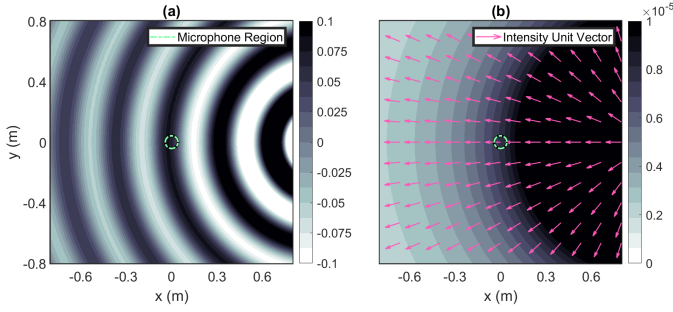


Fig. 4: (a) True pressure field and (b) intensity field at 1000 Hz in the XY-plane with the point-source at $(1, 0, 0)$ m, where intensity magnitude is given by the color-map.

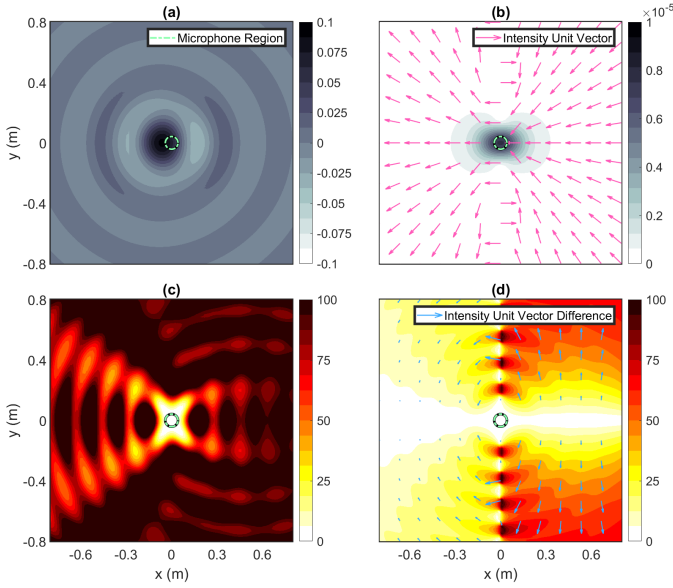


Fig. 5: Truncated measurement of (a) the pressure field and (b) the intensity field at 1000 Hz in the XY-plane for the point-source at $(1, 0, 0)$ m, where (c) is PE and (d) is IDE.

where $I = \frac{1}{2} \text{Re}(PV^*)$, and V^* is the conjugated sound field velocity. The intensity direction error (IDE), which is denoted as the acute angle between the true recorded and reproduced intensity fields [61], is expressed as

$$\text{IDE} = \arccos \left(\frac{I \cdot \tilde{I}}{\|I\| \cdot \|\tilde{I}\|} \right) / \pi \times 100(\%). \quad (37)$$

Additionally, for intensity fields, we also illustrate the true and reproduced intensity unit vector difference $= I/\|I\| - \tilde{I}/\|\tilde{I}\|$.

B. Pressure and Intensity Fields

Figure 4 shows the pressure and intensity field for the true sound-source at $(1, 0, 0)$ m that we recorded and reproduced virtually in the perception experiment. The 4th order recording of this true sound-source is shown in Fig. 5. Immediately we observe the effects of truncation in the recorded pressure field (Fig. 5a), where a distinct near-field pattern is no longer visible. As expected, the recording is seen to be localized

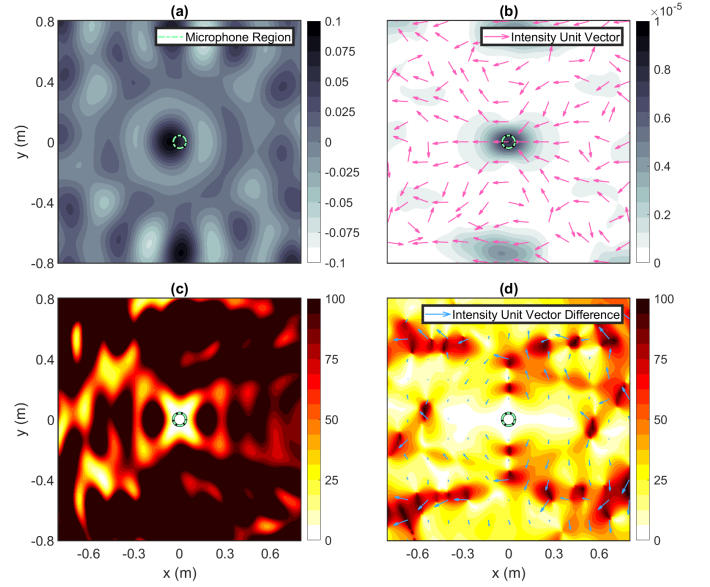


Fig. 6: PW-CF reproduction of (a) pressure field and (b) intensity field at 1000 Hz in the XY-plane, where (c) is the reproduction PE and (d) is the reproduction IDE.

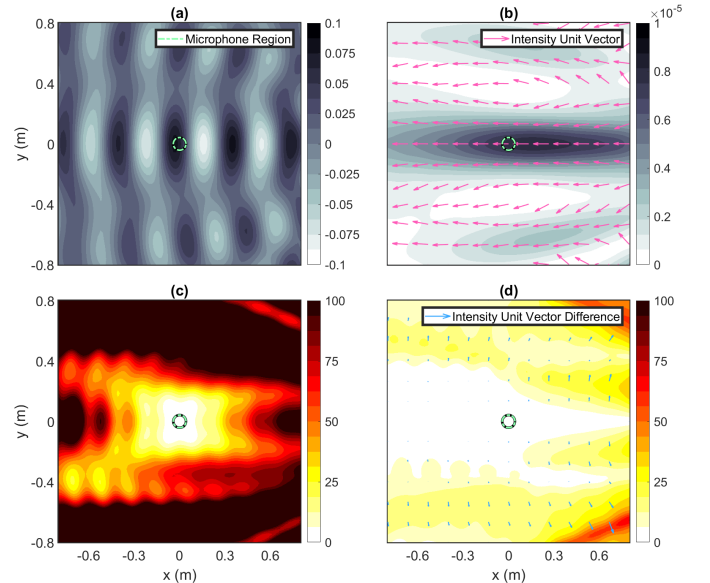


Fig. 7: MW-IRLS reproduction of (a) pressure field and (b) intensity field at 1000 Hz in the XY-plane, where (c) is the reproduction PE and (d) is the reproduction IDE.

spatially within the microphone array, illustrated by the sweet-spot within the PE (Fig. 5c). Similarly, the recorded intensity is seen to be concentrated about the sweet-spot. Beyond the sweet-spot, truncation error is seen to degrade the pressure and intensity accuracy, leading to the perceptual artifacts we wish to resolve by extrapolating a virtual source environment.

In Fig. 6, we observe that the PW-CF experiences the same sweet-spot behaviors as the truncated recording, where once again the reproduced pressure and intensity is localized to the microphone's region (Fig. 6c). A similar result is also obtained by the MW-CF method (not shown), supporting that

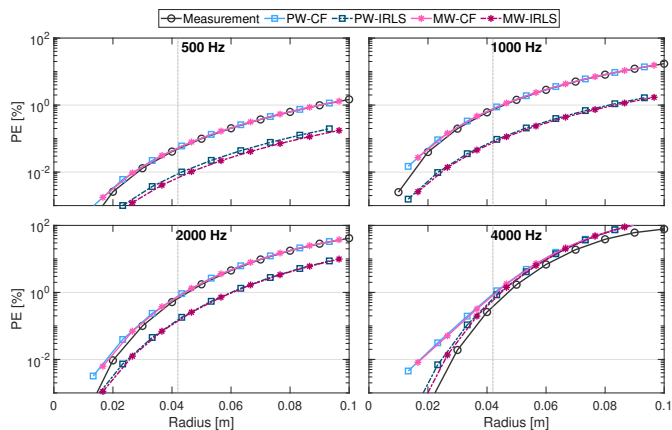


Fig. 8: Average pressure error over a spherical surface of varying radius at four frequencies for the measured and reproduced sound fields.

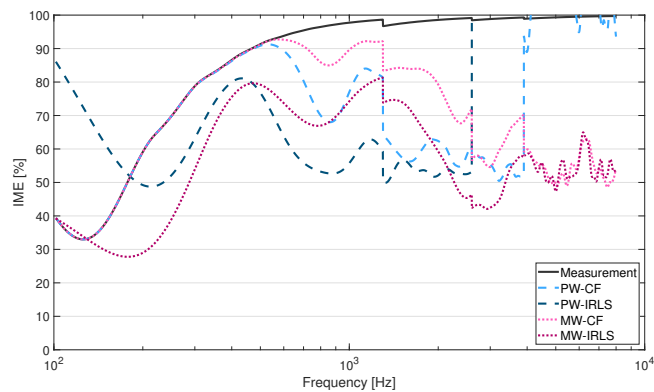


Fig. 9: Average intensity magnitude error over a 0.8 m spherical surface plotted against frequency.

the sweet-spot is caused by the closed-form expansion over-approximating the truncated recording. The PW-CF intensity field is also seen to be non-uniform throughout the virtual environment. It is expected that this may be a dominant factor contributing to the PW-CF’s perceptual evaluation.

Figure 7 shows better results for the MW-IRLS reproduction. As intended, IRLS expansion is seen to relax the sweet-spot constraint (Fig. 7a & c). Similar results are also observed for the PW-IRLS method (not shown), indicating that sparse expansions are able to extend the region of reproduction accuracy. This is believed to aid the perceptual stability of the reproduction as the listener translates further from the original recording position. Furthermore, the MW-IRLS intensity field (Fig. 7b) is shown to have improved uniformity, which leads to better IDE results (Fig. 7d). This uniformity is expected to have contributed to the strong perceptual results achieved by the MW-IRLS method.

C. Pressure and Intensity Error

We present the averaged PE at various translation distances in Fig. 8. A clear difference in performance is observed at the lower frequencies, where the two IRLS expansions (PW-IRLS

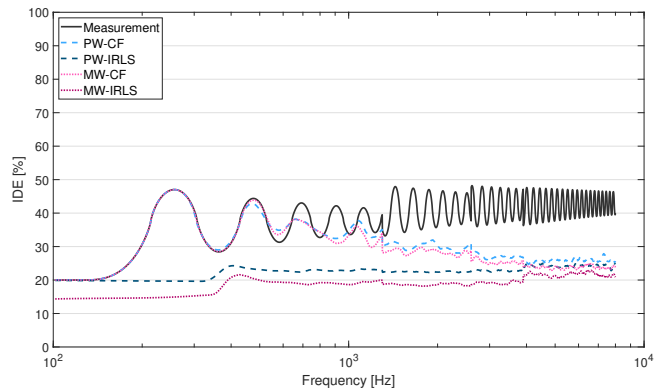


Fig. 10: Average intensity direction error over a 0.8 m spherical surface plotted against frequency.

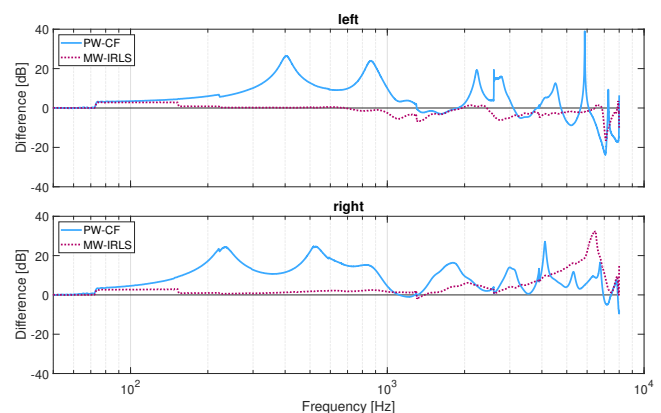


Fig. 11: BRIR spectral difference between the true (reference) and reproduced (PW-CF, MW-IRLS) signals rendered at the translated position $(0, 0.5, 0)$ m.

and MW-IRLS) are seen to better reproduce the pressure field throughout a 0.1 m region. This result corroborates with the prior sweet-spot observations, where the IRLS expansions are able to relax spatial constraints. On the other hand, the closed-form expansions are shown to match the PE of the recording, further illustrating that the PW-CF and MW-CF methods over-approximate the truncation artifacts.

All methods are observed to have poor IME at the translation of 0.8 m in Fig. 9. At higher frequencies both MW-CF and MW-IRLS have lower error than their planewave counterparts. However, the IME is still poor, and it is difficult to know if this behavior contributed to perceptual results. Additionally, large spikes in error are found when the microphone’s truncation increases between the $[k|x_Q|]$ frequency bands. It may be possible to smooth the activation of each band to further improve perceptual stability.

The IDE shows clearer results at the 0.8 m translation in Fig. 10. The MW-IRLS reproduction is seen to strongly match the direction of the true sound-source’s intensity across the full frequency range. This intensity alignment is expected to have contributed to the perceptual results of the MW-IRLS method. This is in contrast to the PW-CF benchmark which is seen to follow the recording’s poor IDE at lower frequencies.

D. BRIR Response

We measure the reproduction system's response by recording, expanding, and auralizing a sine-sweep signal with the planewave and mixedwave translation methods. This provides the binaural room impulse response (BRIR) of the translated listener in the virtual environment. Figure 11 gives the BRIR spectra difference of the PW-CF and MW-IRLS compared to the reference at the translated position of $(0, 0.5, 0)$ m to the left. The BRIR spectral results show the most substantial difference between the PW-CF and MW-IRLS methods thus far. Below 1000 Hz the MW-IRLS is observed to have little spectral deviation from the reference BRIR. This suggests that the MW-IRLS accurately reconstructs the sound heard by a translated listener in the true environment. The PW-CF BRIR, however, is seen to deviate significantly from the reference. Therefore, the BRIR spectra differences support that the MW-IRLS offers greater perceptual accuracy, which was observed in the perceptual experiment results.

The BRIR includes the effects of HRTF processing that is applied to each sound field translation method. This may explain why there is such a large disparity between the planewave and mixedwave BRIR. The mixedwave method adapts the HRTF with the listener's movement. While the planewave method maintains a constant HRTF perspective and shifts the truncated reproduction about the listener. It is the difference between these two HRTF implementations that may have the most significant effect on the BRIR differences and the perceptual experiment results.

VI. CONCLUSION

Virtual reality technology enhances acoustic real-world reproductions by allowing listeners to perceptually move about the environment. At this time, however, the benchmark planewave method towards sound field translation is still limited by inherited microphone constraints. Furthermore, the planewave source model is restricted to the far-field, which results in the listener's HRTF perceptive being fixed during translation. As a result, immersion in the planewave environment is degraded by poor source localizability and audible spectral distortions.

We have proposed an alternative source model for translation that enables a sparse virtual environment to contain a mixture of near-field and far-field sources. We compared this proposed mixedwave method against the planewave benchmark through a perceptual MUSHRA experiment and cross-examined the results with numerical simulations. For human speech reproduction, the mixedwave source model improved both localizability and audio quality. Both the closed-form and IRLS expanded mixedwave reproductions were found to provide a more immersive experience. Similar results were also found for a wider band music sound source.

The IRLS expansion was shown to help enlarge the reproduction sweet-spot, and in response it scored better perceptually than the closed-form expansion for speech sound. Additionally, the proposed method showed better intensity direction matching than the benchmark, further corroborating the perceptual results. Finally, we illustrated that the mixedwave's

ambisonic-like binaural rendering allows for greater perceptual accuracy due to lower BRIR spectral error.

We note that this paper focuses on the sole effects of modeling a near-field far-field mixture for translation. As such, we studied an over-simplified acoustic environment in order to make clearer comparisons. We leave the considerations involved with implementing the proposed method as future work. Accounting for acoustic reflections, diffuse sound, multiple sound sources, source directivity, and the methods to separate and process them in a virtual mixedwave environment are left as an open problem. Furthermore, it may be satisfying to reproduce the recorded experience along with synthetic sounds, and it should be explored along side future applications developed for virtual reality.

VII. THANKS

The authors would like to thank Zamir Ben-Hur for guidance in developing the perceptual test, and Shawn Featherly for development of the perceptual test Unity application.

REFERENCES

- [1] P. Dodds, S. Amengual Garí, W. Brimijoin, and P. Robinson, "Auralization systems for simulation of augmented reality experiences in virtual environments," in *Audio for Virtual, Augmented and Mixed Realities: Proc. ICSA 2019; 5th Intl. Conf. on Spatial Audio*, 2019, pp. 29–34.
- [2] Y. Suzuki *et al.*, "3d spatial sound systems compatible with human's active listening to realize rich high-level kansei information," *Interdisciplinary information sciences*, vol. 18, no. 2, pp. 71–82, 2012.
- [3] J. G. Tylka and E. Y. Choueiri, "Models for evaluating navigational techniques for higher-order ambisonics," in *Proc. Meetings on Acoust. ASA*, 2017, vol. 30, p. 050009.
- [4] S. Amengual Garí, C. Schissler, R. Mehra, S. Featherly, and P. Robinson, "Evaluation of real-time sound propagation engines in a virtual reality framework," in *Proc. Intl. Audio Eng. Soc. Conf. on Immersive and Interactive Audio*. Audio Engineering Society, 2019.
- [5] M. Ziegler *et al.*, "Immersive virtual reality for live-action video using camera arrays," *IBC, Amsterdam, Netherlands*, 2017.
- [6] D. Rivas Méndez, C. Armstrong, J. Stubbs, M. Stiles, and G. Kearney, "Practical recording techniques for music production with six-degrees of freedom virtual reality," in *Audio Eng. Soc. Conv. 145*. Audio Engineering Society, 2018.
- [7] C. D. Salvador, S. Sakamoto, J. Trevino, and Y. Suzuki, "Spatial accuracy of binaural synthesis from rigid spherical microphone array recordings," *Acoust. Sci. Technol.*, vol. 38, no. 1, pp. 23–30, 2017.
- [8] J. G. Tylka and E. Y. Choueiri, "Fundamentals of a parametric method for virtual navigation within an array of ambisonics microphones," *J. Audio Eng. Soc.*, vol. 68, no. 3, pp. 120–137, 2020.
- [9] J. G. Tylka and E. Y. Choueiri, "Performance of linear extrapolation methods for virtual sound field navigation," *J. Audio Eng. Soc.*, vol. 68, no. 3, pp. 138–156, 2020.
- [10] N. Mariette and B. Katz, "Sounddeltalarge scale multi-user audio augmented reality," in *Proc. of the EAA Symposium on Auralization*, 2009, pp. 15–17.
- [11] J. G. Tylka and E. Y. Choueiri, "Domains of practical applicability for parametric interpolation methods for virtual sound field navigation," *J. Audio Eng. Soc.*, vol. 67, no. 11, pp. 882–893, 2019.
- [12] E. Patricio, A. Ruminski, A. Kuklasinski, L. Januszkiwicz, and T. Zernicki, "Toward six degrees of freedom audio recording and playback using multiple ambisonics sound fields," in *Audio Eng. Soc. Conv. 146*. Audio Engineering Society, 2019.
- [13] P. Samarasinghe, T. Abhayapala, and M. Poletti, "Wavefield analysis over large areas using distributed higher order microphones," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 22, no. 3, pp. 647–658, 2014.
- [14] J. G. Tylka and E. Choueiri, "Soundfield navigation using an array of higher-order ambisonics microphones," in *Proc. Intl. Audio Eng. Soc. Conf. on Audio for Virtual and Augmented Reality*. Audio Engineering Society, 2016.

- [15] Y. Wang and K. Chen, "Translations of spherical harmonics expansion coefficients for a sound field using plane wave expansions," *J. Acoust. Soc. Amer.*, vol. 143, no. 6, pp. 3474–3478, 2018.
- [16] O. Thiergart, G. Del Galdo, M. Taseska, and E. Habets, "Geometry-based spatial sound acquisition using distributed microphone arrays," *IEEE Trans. Audio, Speech, Language Process.*, vol. 21, no. 12, pp. 2583–2594, 2013.
- [17] J. G. Tylka and E. Choueiri, "Comparison of techniques for binaural navigation of higher-order ambisonic soundfields," in *Audio Eng. Soc. Conv. 139*. Audio Engineering Society, 2015.
- [18] M. Noisternig, A. Sontacchi, T. Musil, and R. Holdrich, "A 3d ambisonic based binaural sound reproduction system," in *Proc. of 24th Intl. Audio Eng. Soc. Conf. on Multichannel Audio, The New Reality*. Audio Engineering Society, 2003.
- [19] D. Menzies and M. Al-Akaidi, "Ambisonic synthesis of complex sources," *J. Audio Eng. Soc.*, vol. 55, no. 10, pp. 864–876, 2007.
- [20] T. Pihlajamaki and V. Pulkki, "Synthesis of complex sound scenes with transformation of recorded spatial sound in virtual reality," *J. Audio Eng. Soc.*, vol. 63, no. 7/8, pp. 542–551, 2015.
- [21] E. Fernandez-Grande, "Sound field reconstruction using a spherical microphone array," *J. Acoust. Soc. Amer.*, vol. 139, no. 3, pp. 1168–1178, 2016.
- [22] F. Schultz and S. Spors, "Data-based binaural synthesis including rotational and translatory head-movements," in *Proc. of 52nd Intl. Audio Eng. Soc. Conf. on Sound Field Control-Engineering and Perception*. Audio Engineering Society, 2013.
- [23] R. Duraiswami, Z. Li, D. N. Zotkin, E. Grassi, and Nail A. G., "Plane-wave decomposition analysis for spherical microphone arrays," in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust.* IEEE, 2005, pp. 150–153.
- [24] M. A. Poletti, "Three-dimensional surround sound systems based on spherical harmonics," *J. Audio Eng. Soc.*, vol. 53, no. 11, pp. 1004–1025, 2005.
- [25] D. B. Ward and T. D. Abhayapala, "Reproduction of a plane-wave sound field using an array of loudspeakers," *IEEE Trans. Speech Audio Process.*, vol. 9, no. 6, pp. 697–707, 2001.
- [26] MH Acoustics, "Em32 eigenmike microphone array release notes (v17.0)," 25 Summit Ave, Summit, NJ 07901, USA, 2013.
- [27] N. Hahn and S. Spors, "Modal bandwidth reduction in data-based binaural synthesis including translatory head-movements," in *Proc. German Annu. Conf. Acoust.(DAGA)*, 2015, pp. 1122–1125.
- [28] N. Hahn and S. Spors, "Physical properties of modal beamforming in the context of data-based sound reproduction," in *Audio Eng. Soc. Conv. 139*. Audio Engineering Society, 2015.
- [29] A. Kuntz and R. Rabenstein, "Limitations in the extrapolation of wave fields from circular measurements," in *Eur. Signal Process. Conf.* IEEE, 2007, pp. 2331–2335.
- [30] F. Winter, F. Schultz, and S. Spors, "Localization properties of data-based binaural synthesis including translatory head-movements," in *Proceedings of the Forum Acusticum, Krakow, Poland*, 2014, vol. 31.
- [31] J. Daniel, "Spatial sound encoding including near field effect: Introducing distance coding filters and a viable, new ambisonic format," in *Proc. of 23rd Intl. Audio Eng. Soc. Conf. on Signal Processing in Audio Recording and Reproduction*. Audio Engineering Society, 2003.
- [32] K. Wakayama, J. Trevino, H. Takada, S. Sakamoto, and Y. Suzuki, "Extended sound field recording using position information of directional sound sources," in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust.* IEEE, 2017, pp. 185–189.
- [33] A. Plinge, S. Schlecht, O. Thiergart, T. Robotham, O. Rummukainen, and E. Habets, "Six-degrees-of-freedom binaural audio reproduction of first-order ambisonics with distance information," in *Proc. Intl. Audio Eng. Soc. Conf. on Audio for Virtual and Augmented Reality*. Audio Engineering Society, 2018.
- [34] M. Kentgens, A. Behler, and P. Jax, "Translation of a higher order ambisonics sound scene based on parametric decomposition," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2020, pp. 151–155.
- [35] D. Menzies and M. Al-Akaidi, "Nearfield binaural synthesis and ambisonics," *J. Acoust. Soc. Amer.*, vol. 121, no. 3, pp. 1559–1563, 2007.
- [36] Zylia Sp. z o.o., "ZYLIA ZM-1 Microphone," <https://www.zylia.co/zylia-zm-1-microphone.html>, accessed: Feb. 2020.
- [37] VisiSonics Corporation, "VisiSonics 5/64 audio/visual camera," <https://visisonics.com/564avcamera/>, accessed: Feb. 2020.
- [38] R. Chartrand and W. Yin, "Iteratively reweighted algorithms for compressive sensing," in *IEEE Intl. Conf. Acoust., Speech, Signal Process.* IEEE, 2008, pp. 3869–3872.
- [39] S. Emura, "Sound field estimation using two spherical microphone arrays," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* IEEE, 2017, pp. 101–105.
- [40] Y. Hu, P. N. Samarasinghe, T. D. Abhayapala, and G. Dickins, "Modeling characteristics of real loudspeakers using various acoustic models: Modal-domain approaches," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* IEEE, 2019, pp. 561–565.
- [41] Y. Maeno, Y. Mitsufuji, and T. D. Abhayapala, "Mode domain spatial active noise control using sparse signal representation," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* IEEE, 2018, pp. 211–215.
- [42] L. Birnie, T. Abhayapala, P. Samarasinghe, and V. Tourbabin, "Sound field translation methods for binaural reproduction," in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust.* IEEE, 2019, pp. 140–144.
- [43] ITU Radiocommunication Assembly, "ITU-R BS.1534-3: Method for the subjective assessment of intermediate quality level of audio systems," October 2015.
- [44] Z. Ben-Hur, D. L. Alon, B. Rafaely, and R. Mehra, "Loudness stability of binaural sound with spherical harmonic representation of sparse head-related transfer functions," *EURASIP J. Audio, Speech, Music Process.*, vol. 2019, no. 1, pp. 5, 2019.
- [45] E.G. Williams, *Fourier Acoustics: Sound Radiation and Nearfield Acoustic Holography*, Academic Press, London, UK, 1999.
- [46] H. Chen, T. D. Abhayapala, and W. Zhang, "Theory and design of compact hybrid microphone arrays on two-dimensional planes for three-dimensional soundfield analysis," *J. Acoust. Soc. Amer.*, vol. 138, no. 5, pp. 3081–3092, 2015.
- [47] VisiSonics Corporation, "VisiSonics audio/visual planar array," <https://visisonics.com/audio-visual-planar-array/>, accessed: Feb. 2020.
- [48] T. D. Abhayapala and D. B. Ward, "Theory and design of high order sound field microphones using spherical microphone array," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* IEEE, 2002, vol. 2, pp. II–1949.
- [49] B. Rafaely, "Analysis and design of spherical microphone arrays," *IEEE Trans. Speech Audio Process.*, vol. 13, no. 1, pp. 135–143, 2005.
- [50] R. A. Kennedy, P. Sadeghi, T. D. Abhayapala, and H. M. Jones, "Intrinsic limits of dimensionality and richness in random multipath fields," *IEEE Trans. Signal Process.*, vol. 55, no. 6, pp. 2542–2556, 2007.
- [51] D. N. Zotkin, R. Duraiswami, and N. A. Gumerov, "Regularized hrtf fitting using spherical harmonics," in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust.* IEEE, 2009, pp. 257–260.
- [52] W. Zhang, T. D. Abhayapala, R. A. Kennedy, and R. Duraiswami, "Insights into head-related transfer function: Spatial dimensionality and continuous representation," *J. Acoust. Soc. Amer.*, vol. 127, no. 4, pp. 2347–2357, 2010.
- [53] G. N. Lilis, D. Angelosante, and G. B. Giannakis, "Sound field reproduction using the lasso," *IEEE Trans. Audio, Speech, Language Process.*, vol. 18, no. 8, pp. 1902–1912, 2010.
- [54] R. Tibshirani, "Regression shrinkage and selection via the lasso," *J. Roy. Stat. Soc. Ser. B*, vol. 58, no. 1, pp. 267–288, 1996.
- [55] E. J. Candès and M. B. Wakin, "An introduction to compressive sampling," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 21–30, 2008.
- [56] A. Lindau, T. Hohn, and S. Weinzierl, "Binaural resynthesis for comparative studies of acoustical environments," in *Audio Eng. Soc. Conv. 122*. Audio Engineering Society, 2007.
- [57] F. Brinkmann *et al.*, "A cross-evaluated database of measured and simulated hrtfs including 3d head meshes, anthropometric features, and headphone impulse responses," *J. Audio Eng. Soc.*, vol. 67, no. 9, pp. 705–718, 2019.
- [58] F. Brinkmann *et al.*, "The hutubs head-related transfer function (hrtf) database," [online]. <http://dx.doi.org/10.14279/depositonce-8487>, accessed: Feb. 2020.
- [59] B. Rafaely and M. Kleider, "Spherical microphone array beam steering using wigner-d weighting," *IEEE Signal Process. Lett.*, vol. 15, pp. 417–420, 2008.
- [60] J. Fliege and U. Maier, "The distribution of points on the sphere and corresponding cubature formulae," *IMA J. Numer. Anal.*, vol. 19, no. 2, pp. 317–334, 1999.
- [61] M. Shin, P. A. Nelson, F. M. Fazi, and J. Seo, "Velocity controlled sound field reproduction by non-uniformly spaced loudspeakers," *Journal of Sound and Vibration*, vol. 370, pp. 444–464, 2016.