

Time-difference of Arrival Model for Spherical Microphone Arrays and Application to Direction of Arrival Estimation

Joonas Nikunen, and Tuomas Virtanen,

Laboratory of Signal Processing, Tampere University of Technology, Tampere, Finland
joonas.nikunen@tut.fi, tuomas.virtanen@tut.fi

Abstract—This paper investigates different steering techniques for spherical microphone arrays and proposes a time-difference of arrival (TDOA) model for microphones on surface of a rigid sphere. The model is based on geometric interpretation of wavefront incident angle and the extra distance the wavefront needs to travel to reach microphones on the opposite side of a sphere. We evaluate the proposed model by comparing analytic TDOAs to measured TDOAs extracted from impulse responses (IR) of a rigid sphere ($r = 7.5\text{cm}$). The proposed method achieves over 40% relative improvement in TDOA accuracy in comparison to free-field propagation and TDOAs extracted from analytic IRs of a spherical microphone array provide an additional 10% improvement. We test the proposed model for the application of source direction of arrival (DOA) estimation using steered response power (SRP) with real reverberant recordings of moving speech sources. All tested methods perform equally well in noise-free scenario, while the proposed model and simulated IRs improve over free-field assumption in low SNR conditions. The proposed model has the benefit of only using single delay for steering the array.

I. INTRODUCTION

Direction of arrival (DOA) estimation of a sound source using microphone array is based on spatiotemporal filtering methods such as steered response power (SRP) [1]–[3]. Processing with spaced arrays assumes direct-path propagation from sources to microphones, however, if the array consists of microphones embedded on a surface of a rigid sphere the free-field assumption is violated. The analytic properties of sound wave scattering on a surface of rigid sphere are well known [4], [5] and can be utilized by processing using the spherical harmonic decomposition [6]. However, it is not deemed feasible with spherical arrays embodying only a small amount of capsules leading to uneven and sparse distribution of microphones.

The properties of spherical arrays can be also considered in the space domain by using analytic spherical array impulse responses (IR) or by simplification down to a single delay caused by the wavefront curvature around the sphere. This can be achieved by measuring IRs for a given array towards all directions and extracting the realized delay from the direct-path peaks of the IRs. However, this is very labor and time consuming if a high number of directions is used and the

process needs to be repeated if the microphone placement or the radius of the sphere is changed. Alternatively, solving the analytic IRs for microphones on the surface of a sphere [4], [5] allows using the obtained IRs for steering the array by convolution.

For the delay based processing we propose a time-difference of arrival (TDOA) model for microphones on a surface of a rigid sphere. The proposed model is based on the distance and corresponding travel time the wavefront is required to curve in order to reach the microphones on the opposite side of the sphere. The applications of the proposed method include all TDOA-based spatial audio signal processing, e.g., DOA estimation [7], source tracking [8] and source separation [9], [10]. The benefit of proposed TDOA model in comparison to using analytic IRs for steering the array is the low computational complexity in time-domain processing where the steering could be done by a single delay element. For example, localization of sound sources with miniature robots require very low computational complexity due to their limited CPU performance and required energy efficiency. However, in the frequency domain the computational cost of using simulated IRs or single time delay are the same, since the required convolution is carried out by multiplication between array signals and frequency domain steering vector.

The evaluation is based on comparing the proposed analytic TDOAs to actual measured TDOAs of spherical array embodying eight microphones. A 40% relative improvement in the TDOA accuracy is achieved with the proposed model in comparison to free-field assumption. The TDOAs extracted from simulated spherical array IRs improve the accuracy by an additional 10%. We demonstrate the performance of the proposed model for source DOA estimation with real recordings of moving speech sources. The DOA estimation accuracy of the proposed model show improved performance in noisy conditions over free-field assumption and it achieves similar performance in comparison to simulated spherical array IRs.

The rest of the paper is organized as follows. First we introduce the proposed TDOA model for microphones on a surface of a rigid sphere in Section II. The evaluation of error between measured TDOA and the different analytic TDOAs is presented in Section III. The impact of the proposed model for the application of source DOA estimation using SRP is

This research was supported by Nokia Technologies.

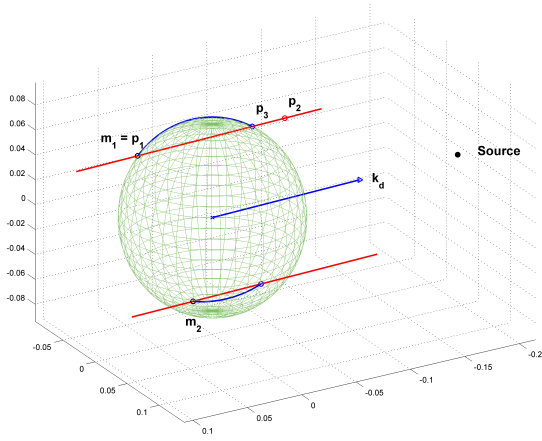


Fig. 1. Sphere-line intersection: Surface of a rigid sphere depicted using mesh, black circles mark the two microphones, blue arrow denotes direction vector \mathbf{k}_d , points of sphere-line intersection denoted by blue circles and the great circle distances illustrated by blue arcs.

evaluated in Section IV. The work is concluded in Section V.

II. TDOA MODEL FOR MICROPHONES ON A RIGID SPHERE

A. Time-difference of arrival in free-field

Spatial audio signal processing is based on observing time and level difference between the microphone elements. A sound wave emitted by a sound source can be considered as a plane wave (if the source is far enough from the array), and the source direction corresponds to a set of time-delays between microphones determined by the array geometry (location of microphones).

In this paper the DOA of the wavefront is denoted using a unit direction vector $\mathbf{k}_d \in \mathcal{R}^3$, $\|\mathbf{k}_d\| = 1$ which points towards a source at a direction indexed by d parametrized by azimuth $\theta_d \in [0, 2\pi]$ and elevation $\varphi_d \in [0, \pi]$. The direction vector originates from the geometric center of the array \mathbf{p} , which is set to be in the origin ($\mathbf{p} = [0, 0, 0]^T$) in order to simplify the equations. We define a microphone array consisting of two microphones m_1 and m_2 at locations $\mathbf{m}_1 \in \mathcal{R}^3$, $\mathbf{m}_2 \in \mathcal{R}^3$, respectively. Given a source at direction \mathbf{k}_d and assuming the free-field propagation, the TDOA between the microphone pair (m_1, m_2) equals to

$$\tau_d(m_1, m_2) = \frac{-\mathbf{k}_d \cdot (\mathbf{m}_2 - \mathbf{m}_1)}{v}, \quad (1)$$

where v is the speed of sound and \cdot denotes inner product. The above formulation can be interpreted as Euclidean distance between the microphones projected on the direction vector.

B. Time-difference of arrival on a surface of a rigid sphere

The geometric interpretation of proposed model equals to well known Woodworth model [11], [12] for computing interaural time-difference between ears by ray-tracing and considering human head as a rigid sphere. Assuming that a microphone on a surface of a sphere is located in such way that it has no direct line of sight for the source denoted by the direction vector. The proposed method is based on analytically

finding the closest point with respect to the microphone that the plane wave can reach in free-field propagation (i.e. without curvature) and accounting for the additional distance it is required to travel on the surface of the sphere to reach the microphone. This equals to the great circle distance (GCD), which denotes shortest distance between two points on the surface of a sphere. The additional distance and travel time corresponding to any microphone location and direction vector can be obtained by solving the points of line-sphere intersection, where the line is parallel to the direction vector and pass through the location of the microphone. The wavefront curvature happens exactly at the halfway of the GCD connecting the microphone location and point of intersection.

1) *Line-sphere intersection*: Two microphones on a surface of a rigid sphere and a direction vector \mathbf{k}_d are depicted in Figure 1. In the following, the method of obtaining the proposed travel time correction is presented for a single microphone (\mathbf{m}_1). The process is repeated similarly for the second microphone (\mathbf{m}_2). We first define a line that passes through the location of microphone \mathbf{m}_1 and direction of the line is parallel to the look direction vector \mathbf{k}_d . The line is depicted in red in Figure 1. The intersecting line can be defined using two points. The first point is the location of the microphone $\mathbf{p}_1 = \mathbf{m}_1$ and the second is obtained as $\mathbf{p}_2 = \mathbf{m}_1 + \mathbf{k}_d$. In order to simplify equations we omit the direction index d for the time being and the resulting line can be parametrized as

$$\mathbf{y} = \mathbf{o} + l\mathbf{u}, \quad (2)$$

where $\mathbf{u} = \frac{\mathbf{p}_2 - \mathbf{p}_1}{\|\mathbf{p}_2 - \mathbf{p}_1\|}$ is the unit direction vector and \mathbf{y} are the points in line parametrized by the distance l from the line origin \mathbf{o} . The sphere can be expressed as

$$\|\mathbf{y} - \mathbf{c}\|^2 = r^2 \quad (3)$$

where \mathbf{c} is the center of the sphere and r is its radius. One can easily solve the points of line-sphere intersection as

$$l = -(\mathbf{u} \cdot (\mathbf{o} - \mathbf{c})) \pm \sqrt{(\mathbf{u} \cdot (\mathbf{o} - \mathbf{c}))^2 - \|\mathbf{o} - \mathbf{c}\|^2 + r^2}. \quad (4)$$

2) *Proposed TDOA model*: In order to avoid for searching if the microphone is at the same hemisphere as the source and no correction for travel time is needed (direct line of sight), we now set \mathbf{p}_2 as the origin of the coordinate system. Solving for l with the new origin using (4) and finding the intersection point closest to origin (smaller absolute value of the two solutions, denoted hereafter as \hat{l}) equals to obtaining the desired intersection point \mathbf{p}_3 if the microphone is truly in the opposite side of sphere. In case of microphone being on the same hemisphere as the source we obtain the location of the microphone as the intersection point, which further indicates that no correction is needed. The resulting intersection point can be given as

$$\mathbf{p}_3 = \mathbf{p}_2 + \hat{l}\mathbf{u}, \quad (5)$$

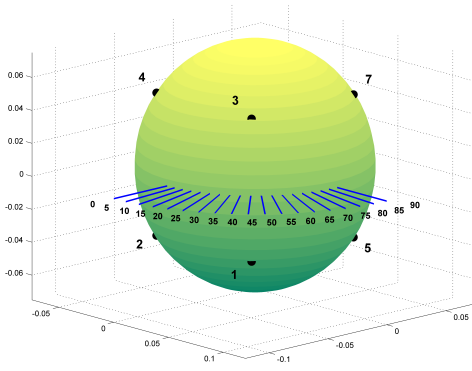


Fig. 2. Illustration of the spherical microphone array used and the directions from where the IRs were measured from. Actual array was 3D printed with holes for miniature microphone capsules.

TABLE I
MEAN ABSOLUTE TDOA ERROR IN SAMPLES AT $F_s = 48$ KHZ.

	All pairs	Mic. 6 omitted	Non-direct pairs
FF	0.939 (100%)	0.614 (100%)	1.122 (100%)
GCD	0.545 (58%)	0.343 (56%)	0.622 (55%)
SPH	0.402 (43%)	0.291 (47%)	0.447 (40%)

depicted in the Figure 1 by a blue circle. Having obtained the intersection point \mathbf{p}_3 one can easily calculate the GCD between it and the microphone location $\mathbf{p}_1 = \mathbf{m}_1$ as

$$d_{gcd} = r \arctan \frac{|\mathbf{p}_1 \times \mathbf{p}_3|}{\mathbf{p}_1 \cdot \mathbf{p}_3}. \quad (6)$$

The resulting arc connecting the two points is denoted by blue in the Figure 1.

As pointed out earlier, the point of curvature for the plane wave happens exactly at the halfway of the GCD (d_{gcd}) and thus additional distance is only accounted for the half of d_{gcd} resulting to TDOA correction of

$$\hat{\tau} = \frac{d_{gcd}/2 - d_{Euc}/2}{v}, \quad (7)$$

where d_{Euc} denotes Euclidean distance between points \mathbf{p}_1 and \mathbf{p}_3 . Reintroducing the direction index d , the TDOA correction towards direction \mathbf{k}_d for a microphone m_1 is denoted as $\hat{\tau}_d(m_1)$. The resulting proposed TDOA model is defined as

$$\tau_d^{gcd}(m_1, m_2) = \frac{-\mathbf{k} \cdot (\mathbf{m}_2 - \mathbf{m}_1)}{v} - \hat{\tau}_d(m_1) + \hat{\tau}_d(m_2). \quad (8)$$

III. EVALUATION WITH MEASURED TDOAS

We first evaluate the proposed TDOA model using ground truth TDOAs extracted from a IR measurements with a 3D printed spherical microphone array ($r = 7.5$ cm) embodying eight DPA 4060 miniature omnidirectional microphones. An illustration of the sphere and microphones at corners of an imaginary cube ($a = \sqrt{(2r)^2/3}$) is given in Figure 2 along with angles indicating the directions from where the IRs were measured from.

The IRs were measured in sound-insulated and acoustically treated room with average reverberation time of $T_{60} = 260$ ms. The moderate reverberation does not affect the TDOA analysis

by cross-correlation, since the direct path peak dominates over reflections and scattering. Azimuth angles from 0 to 90 degrees with 5 degree spacing at zero elevation were measured. The measured IRs correspond to one unique 90-degree segment of the full sphere and other segments are redundant due to the symmetric placement of microphones.

A Genelec G Two loudspeaker was used for playback and IRs were measured using 18-cycle MLS-sequence as measurement signal. The alignment of the loudspeaker with respect to sphere was achieved with a laser pointer. The mismatch between visual alignment and desired angle was compensated by analyzing the DOA from the recorded signals using SRP [13] and manually picking the maximum indicating the realized DOA. The ground truth TDOAs were estimated by finding the maximum cross-correlation of IRs between microphone pairs. In order to increase the ground truth TDOA resolution, oversampling of the IRs was done using cubic interpolation by a factor of four, resulting to 1/4 subsample accuracy. The simulated spherical array IRs were obtained using spherical array processing library implemented and made freely available by author of [14]. The process of obtaining the TDOAs from the simulated IRs is the same as for real IR measurements.

Absolute errors with respect to ground truth TDOA for each direction and each microphone pair using analytic free-field TDOA (1), the proposed model (8) and the simulated spherical array IRs [14] are illustrated in Figure 3 (a-c), respectively. The results are reported by converting the TDOA error to samples at $F_s = 48$ kHz. The analyzed DOA of each IR used in specifying the direction for analytic TDOAs and simulated IRs is reported in the x-axis of the Figure 3 (a-c).

The results indicate significant decrease of TDOA error with the use of proposed model and simulated IRs in the case of angles and microphone pairs including a non-direct microphone (5,6,7,8 with angles $< 45^\circ$ and 2,4,6,8 with angles $> 45^\circ$). Pair-wise results in Figure 3 indicate that the TDOA error is largest with microphone pairs including microphone 6 and a visual inspection of the IRs from it indicated weak direct path peak. We assume that it is a result of defective placement of the microphone and its diaphragm being below the sphere surface resulting to less reliable ground truth TDOA.

The mean absolute TDOA error over all microphone pairs and all locations is reported in Table I and the methods are abbreviated as FF for free-field assumption, GCD for proposed model and SPH for the analytic spherical array IRs. Also results by omitting pairs including microphone 6 and only using non-direct pairs are reported. The relative improvement in mean absolute TDOA error is greater than 40% when using the proposed model in comparison to FF. The TDOAs analyzed from simulated IRs decrease the error by another 10%.

IV. APPLICATION TO DOA ESTIMATION

A. DOA estimation using steered response power

A microphone array with $m = 1, \dots, M$ microphones is used to capture a signal of a moving sound source along

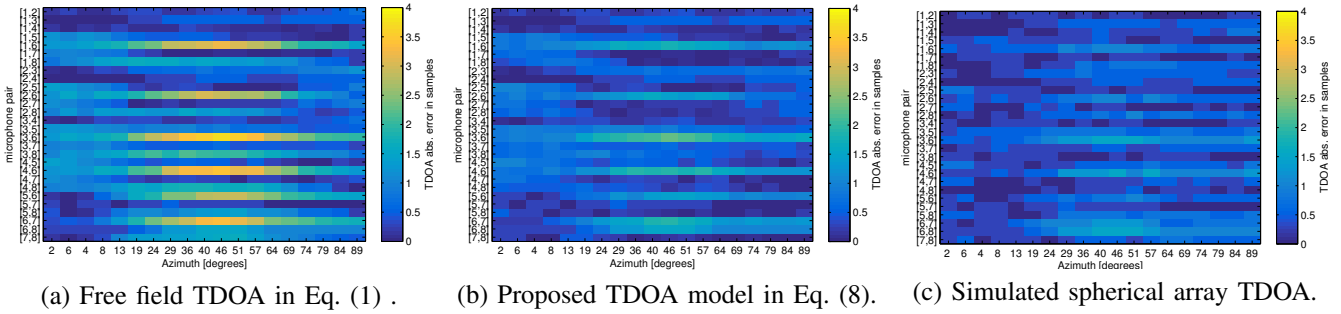


Fig. 3. Absolute TDOA error in samples in comparison to ground truth TDOA measured from IRs.

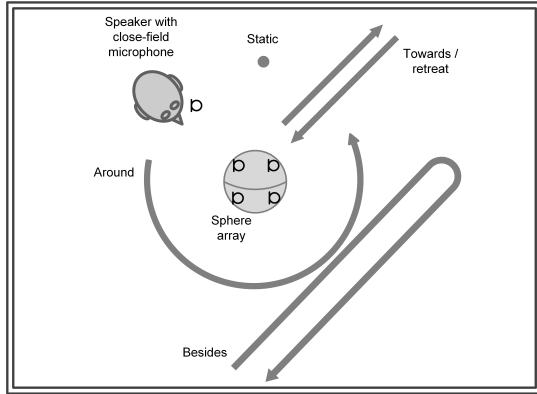


Fig. 4. Recording setup and source movement patterns.

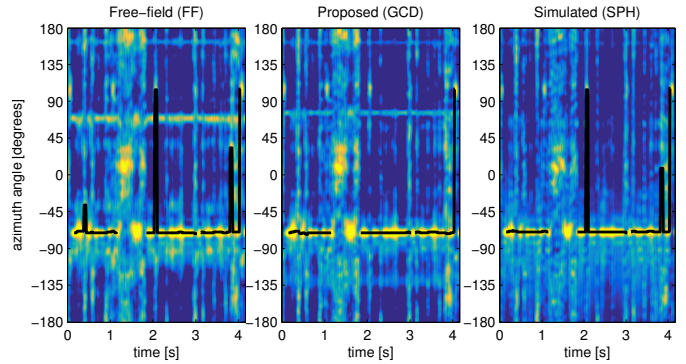


Fig. 5. Example SRP with corresponding DOA trajectories.

with additive uncorrelated noise. The short time Fourier transform (STFT) of the array signal is denoted by $\mathbf{x}_{fn} = [x_{fn1}, \dots, x_{fnM}]^T$ where f and n are frequency and frame index, respectively. Steered response power (SRP) [13] denotes the spatial energy of the captured signal as a function of DOA for each time frame and can be calculated with PHAT weighting as

$$S_{dn} = \sum_{m_1=1}^{M-1} \sum_{m_2=m_1+1}^M \sum_{f=1}^F \frac{x_{fnm_1} x_{fnm_2}^*}{|x_{fnm_1} x_{fnm_2}^*|} \exp(j\omega_f \tau_d(m_1, m_2)), \quad (9)$$

where $*$ denotes complex-conjugate and the term $\exp(j\omega_f \tau_d(m_1, m_2))$ is responsible for time-aligning the microphone signals ($\omega_f = 2\pi(f-1)F_s/N$ where N is the STFT window length). Searching for the maximum of the SRP function at each time frame n corresponds to DOA estimation of the source dominating the spatial evidence in that frame.

B. Task and material collection

We evaluate the performance of DOA estimation by finding maximum of SRP function at each time frame and comparing obtained DOA to the annotated ground truth. The test material was collected with the earlier introduced spherical array ($r = 7.5$ cm) and the place of recording was a typical office building lounge area with moderate reverberation (> 300 ms) and irregular walls and furnishing. Array capture of a speech source moving around the array or being stationary

was recorded and the movement paths are illustrated in Figure 4. The ground truth source DOA trajectories were annotated by hand based on the SRP energy maps calculated using the free-field TDOA. Three different speakers spoke phonetically balanced sentences and in total 12 signals each with 60-second duration were recorded.

For the processing, the recorded signals were downsampled to $F_s = 24$ kHz and the STFT was calculated in windows of 512 samples with 50% frame overlap. The SRP function in Equation (9) was calculated at zero elevation and 1 degree azimuthal spacing between scanned directions. The resulting SRP was averaged over time (9 frames) and after search of maximum index (i.e. the DOA) from each frame the resulting trajectories were median filtered using 4 past and 4 future DOA values. The simulated spherical array IRs and corresponding steering vectors were again obtained using the spherical array processing library provided in [14].

C. DOA estimation results

An example of SRP function by using free-field TDOA (1), proposed TDOA (8) and simulated IRs as steering vectors [14] along with the resulting DOA trajectories is illustrated in Figure 5 (a). One can clearly see the effect of the proposed processing attenuating phantom peaks in the SRP function and the simulated IRs remove the phantom SPR evidence completely around 75° .

Evaluation of the estimated source DOA trajectories is achieved by calculating the mean absolute error (MAE) between the estimated DOA trajectory and the ground truth an-

TABLE II
DOA ESTIMATION RESULTS WITH DIFFERENT CRITERIA.

SNR	MAE [degrees]				SRP peak-%				$\pm 10^\circ$ DOA-%			
	FF	GCD	SPH	SPH ¹	FF	GCD	SPH	SPH ¹	FF	GCD	SPH	SPH ¹
∞	10.7	10.9	9.3	10.5	17.7	19.6	18.7	19.9	88.6	88.4	88.7	88.6
35	10.5	10.8	9.3	10.2	17.3	19.1	18.1	19.4	88.8	88.5	88.8	88.8
30	9.8	10.3	9.2	9.8	16.8	18.5	17.4	18.8	89.3	89.1	89.1	88.8
25	9.5	9.4	9.6	9.6	15.9	17.4	16.3	17.8	89.2	89.5	88.7	88.9
20	9.3	8.9	10.0	8.8	14.9	16.1	14.9	16.4	87.2	88.0	88.0	88.9
15	10.2	9.5	12.6	8.9	13.6	14.5	13.6	14.8	83.6	85.4	82.9	86.0
10	13.0	11.8	18.4	10.7	12.2	12.9	12.2	13.1	77.4	79.4	74.2	80.4
5	19.4	17.5	27.9	16.8	10.9	11.4	11.0	11.5	65.3	68.8	61.9	70.3
0	30.7	28.0	43.1	26.8	9.7	9.9	9.9	10.1	50.0	54.2	44.8	55.5
AVG	13.7	13.0	16.6	12.4	14.3	15.5	14.7	15.8	79.9	81.2	78.6	81.8

notation and averaging over all test signals. The trajectories are only evaluated at frames where speech source is active in order to avoid evaluation from frames consisting only noise. The voice activity detection was achieved by energy thresholding the close-field speech signal recorded during the recording. Also the percentage of correct DOA estimated within $\pm 10^\circ$ from the annotations are reported ($\pm 10^\circ$ DOA-%). It reduces the effect of small errors in the ground truth trajectories by considering estimation being correct within the vicinity of the annotated DOA. Additionally, we measure the SRP evidence mass around the maximum peak (± 10 degrees) with respect to entire SRP evidence integrated over all directions. The **SRP peak-%** measures the peak concentration and amount of phantom spatial evidence at all other directions. Low SRP peak-% indicates high probability of phantom or reflected SRP evidence become chosen as the DOA estimate.

The compared methods are same as earlier with the addition of condition SPH¹ which corresponds to using only the phase of the steering vector obtained from the simulated IRs. The DOA estimation results with varying amount of white Gaussian noise added to the array signals (∞ to 0 dB SNR) are reported in Table II. The proposed model improves DOA estimation accuracy over FF at noisy conditions (SNR < 25 dB). In high SNR conditions FF and GCD have similar DOA estimation performance with the SRP peak-% improved in favor of GCD. The SPH has best MAE performance in noise-free conditions, while with SNRs below 20 dB its performance starts decreasing rapidly. This is due to possible amplification of uncorrelated noise by the simulated IRs whereas FF and GCD only delay the signals before summation in the SRP calculation.

The noise amplification of SPH can be avoided by only using phase of the steering vector (condition SPH¹), which is equivalent to only delaying the signals. However, in comparison to GCD the different frequencies can be delayed by different amounts according to the spherical array theory. The SPH¹ condition results to best overall performance and clearly highest SRP-% with the proposed method right behind it. It can be concluded that the proposed TDOA based steering performs almost as good as the simulated IRs, and for low computational resource applications operating on time-domain

could be achieved without any significant compromise in DOA estimation performance.

V. CONCLUSION

We proposed a TDOA model for spherical microphone array and provided an evaluation of the proposed model by comparison to real measured TDOAs indicating a solid improvement over analytic free-field TDOAs. We also demonstrated the applicability of the proposed model by providing evaluation of its performance in source DOA estimation task. The results indicated improved performance over free-field assumption especially in noisy scenario and similar performance in comparison to steering with simulated spherical array IRs. The benefit of the proposed model is the possible computational efficiency resulting from using only a single delay to align the array signals.

REFERENCES

- [1] Jacek Dmochowski, Jacob Benesty, and Sofine Affes, "Direction of arrival estimation using the parameterized spatial correlation matrix," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 4, pp. 1327–1339, 2007.
- [2] Joseph H DiBiase, Harvey F Silverman, and Michael S Brandstein, "Robust localization in reverberant rooms," in *Microphone Arrays*, Michael S Brandstein and Darren Ward, Eds., pp. 157–180. Springer, 2001.
- [3] Jacek P Dmochowski and Jacob Benesty, "Steered beamforming approaches for acoustic source localization," in *Speech processing in modern communication*, pp. 307–337. Springer, 2010.
- [4] Jens Meyer and Gary Elko, "A highly scalable spherical microphone array based on an orthonormal decomposition of the soundfield," in *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*. IEEE, 2002, vol. 2, pp. 1781–1784.
- [5] Munhum Park and Boaz Rafaely, "Sound-field analysis by plane-wave decomposition using spherical microphone array," *The Journal of the Acoustical Society of America*, vol. 118, no. 5, pp. 3094–3103, 2005.
- [6] Haohai Sun, Heinz Teutsch, Edwin Mabande, and Walter Kellermann, "Robust localization of multiple sources in reverberant environments using eb-esprit with spherical microphone arrays," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2011, pp. 117–120.
- [7] Jacek P Dmochowski, Jacob Benesty, and Sofiene Affes, "A generalized steered response power method for computationally viable source localization," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 8, pp. 2510–2526, 2007.
- [8] Ulrich Klee, Tobias Gehrig, and John McDonough, "Kalman filters for time delay of arrival-based source localization," *EURASIP Journal on Applied Signal Processing*, vol. 2006, pp. 167–167, 2006.
- [9] Joonas Nikunen and Tuomas Virtanen, "Direction of arrival based spatial covariance model for blind sound source separation," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 3, pp. 727–739, 2014.
- [10] Johannes Traa, Paris Smaragdis, Noah D Stein, and David Wingate, "Directional nmf for joint source localization and separation," in *Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2015.
- [11] Robert S Woodworth and Harold Schlosberg, "Experimental psychology (rev)," 1954.
- [12] Richard O Duda and William L Martens, "Range dependence of the response of a spherical head model," *The Journal of the Acoustical Society of America*, vol. 104, no. 5, pp. 3048–3058, 1998.
- [13] I.J. Tashev, *Sound Capture and Processing: Practical Approaches*, John Wiley & Sons Inc, 2009.
- [14] Archontis Politis, *Microphone array processing for parametric spatial audio techniques*, Ph.D. thesis, Aalto University, 2016.