



Aalborg Universitet

AALBORG UNIVERSITY  
DENMARK

## Real-time Loudspeaker Distance Estimation with Stereo Audio

Nielsen, Jesper Kjær; Gaubitch, Nikolay; Heusdens, Richard; Martinez, Jorge; Jensen, Tobias Lindstrøm; Jensen, Søren Holdt

*Published in:*  
23rd European Signal Processing Conference (EUSIPCO), 2015

*DOI (link to publication from Publisher):*  
[10.1109/EUSIPCO.2015.7362383](https://doi.org/10.1109/EUSIPCO.2015.7362383)

*Publication date:*  
2015

*Document Version*  
Accepted author manuscript, peer reviewed version

[Link to publication from Aalborg University](#)

*Citation for published version (APA):*  
Nielsen, J. K., Gaubitch, N., Heusdens, R., Martinez, J., Jensen, T. L., & Jensen, S. H. (2015). Real-time Loudspeaker Distance Estimation with Stereo Audio. In 23rd European Signal Processing Conference (EUSIPCO), 2015 (pp. 250 - 254). [7362383] IEEE Press. Proceedings of the European Signal Processing Conference <https://doi.org/10.1109/EUSIPCO.2015.7362383>

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- ? Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- ? You may not further distribute the material or use it for any profit-making activity or commercial gain
- ? You may freely distribute the URL identifying the publication in the public portal ?

### Take down policy

If you believe that this document breaches copyright please contact us at [vbn@aub.aau.dk](mailto:vbn@aub.aau.dk) providing details, and we will remove access to the work immediately and investigate your claim.

# REAL-TIME LOUDSPEAKER DISTANCE ESTIMATION WITH STEREO AUDIO

*Jesper Kjør Nielsen*<sup>1,2</sup>, *Nikolay D. Gaubitch*<sup>3</sup>, *Richard Heusdens*<sup>3</sup>, *Jorge Martinez*<sup>3</sup>,  
*Tobias Lindstrøm Jensen*<sup>1</sup>, and *Søren Holdt Jensen*<sup>1</sup>

<sup>1</sup>Aalborg University, Denmark  
Dept. of Electronic Systems  
{jkn,tlj,shj}@es.aau.dk

<sup>2</sup>Bang & Olufsen A/S  
Struer, Denmark

<sup>3</sup>Delft University of Technology, The Netherlands  
Signal and Information Processing Lab  
{n.d.gaubitch,r.heusdens,  
j.a.martinezcastaneda}@tudelft.nl

## ABSTRACT

Knowledge on how a number of loudspeakers are positioned relative to a listening position can be used to enhance the listening experience. Usually, these loudspeaker positions are estimated using calibration signals, either audible or psycho-acoustically hidden inside the desired audio signal. In this paper, we propose to use the desired audio signal instead. Specifically, we treat the case of estimating the distance between two loudspeakers playing back a stereo music or speech signal. In this connection, we develop a real-time maximum likelihood estimator and demonstrate that it has a variance in the millimetre range in a real environment for even a modest sampling frequency.

**Index Terms**— Loudspeaker localisation, distance estimation, time-of-arrival estimation

## 1. INTRODUCTION

The distribution of a number of loudspeakers relative to the listening position has a large impact on the listening experience and the perceived spaciousness of sound [1–3]. Often, however, the loudspeakers are not placed in the optimal position since other interior design considerations take higher priority or the desired listening position moves. This can to some extent be compensated for by preprocessing the loudspeaker signals. However, in order to apply the correct preprocessing, the location of the loudspeakers relative to the listening position must be known.

Existing approaches to solving this loudspeaker localisation problem can roughly be dichotomised into two groups. In the first group, synthetic test signals such as sinusoidal sweeps or maximum length sequences (MLS) are used as calibration signals [4–6]. This has the advantage of high estimation accuracy, but also requires the user to actively start the calibration sequence every time, e.g. the listening position or the loudspeaker locations change. This is solved in the second group of methods by adding a calibration signal to the desired audio signal [7, 8]. The calibration signal is shaped psycho-acoustically and hidden inside the audio signal so that it is inaudible to the listener. Consequently, the

energy of the calibration signal is low compared to the energy of the audio signal. This is a problem since the audio signal is considered to be "noise" in the source localisation algorithm, and this affects the estimation accuracy [8].

So why not use the audio signal for source localisation? To our surprise, we have not been able to find any previous work on this. We believe that the main reason for this is that audio signals are much more difficult to work with since they are heavily correlated in both time and in between the loudspeaker channels and have an unknown frequency content. Consequently, it is hard to estimate impulse responses, and the simple cross-correlation methods for loudspeaker localisation fail. Synthetic calibration signals, on the other hand, can be designed to be uncorrelated and to have a desirable frequency content. Thus, the simple cross-correlation methods and impulse response peak picking can be used to compute the distances and/or direction of arrivals (DOAs) between the loudspeakers and/or to the listening position.

In this paper, we take a first step in the direction of loudspeaker localisation using only the desired audio signals. Specifically, we focus on the case where we have to estimate the distance between two loudspeakers playing back a stereo music signal. Distances between all the loudspeaker pairs in a set of loudspeakers can be used to form an Euclidean distance matrix to which the positions of the loudspeakers can be fitted using, e.g. the multidimensional scaling (MDS) algorithm [9] or the algorithm by Crocco et al. [10]. The latter method has the advantage that the loudspeakers and microphones do not have to be co-located. Here, however, we assume that a microphone is mounted on every loudspeaker, which we will refer to as a transceiver, so that they are approximately co-located. This assumption is used in the proposed estimator of the distance to take into account that both transceivers in a transceiver pair should measure the same distance. This increases the robustness of the estimator. Inspired by the recent work in [11], we also formulate the signal model so that the estimator produces estimates from a continuous set without resorting to any heuristic interpolation method. This is in contrast to many of the proposed localisation methods whose resolution is bounded to the sampling grid (see, e.g., [4, 5, 7]).

## 2. ESTIMATING THE DISTANCE BETWEEN TWO TRANSCEIVERS

As alluded to in the introduction, we focus on estimating the distance between two transceivers playing back stereo music or a speech signal. In this paper, a transceiver is a loudspeaker with a microphone mounted close to the diaphragm of the loudspeaker. The developed estimator is not only limited in scope to this special case, but can also be used for the problem where the direct distance should be estimated from a loudspeaker to a microphone, e.g, placed at the listening position, and for the problem where the distance to a reflector should be estimated using just one transceiver. These special cases are obtained by appropriately selecting the source and sensor signals.

### 2.1. The Signal Model

We assume that the two transceivers record  $N$  samples each, and we model these as

$$x_1(n) = q_{11}(n) + q_{21}(n) + e_1(n) \quad (1)$$

$$x_2(n) = q_{22}(n) + q_{12}(n) + e_2(n) \quad (2)$$

where  $e_i(n)$  and  $q_{ki}(n)$  are the noise recorded by transceiver  $i$  and the signal recorded by transceiver  $i$  from transceiver  $k$ , respectively. Thus,  $q_{ii}(n)$  is the part of the microphone signal  $x_i(n)$  which originates from transceiver  $i$ . This signal is not of interest as it does not contain any information on the distance between the transceivers, and we therefore wish to suppress it as much as possible. To do that, we model  $q_{ii}(n)$  as

$$q_{ii}(n) = \sum_{m=0}^{M-1} h_i(m)s_i(n-m) \quad (3)$$

where  $s_i(n)$  and  $h_i(m)$  are a source signal sample of transceiver  $i$  and an FIR filter coefficient of the  $i$ th  $M$ -length transceiver filter, respectively. Thus, a transceiver filter models the acoustic impulse response between the loudspeaker and microphone on a transceiver. We assume that the loudspeakers and microphones are all connected to the same system so that the source signals are known. On the other hand, the transceiver filters are assumed unknown since these might be slowly time-varying due to, e.g., temperature changes. These transceiver filters are very important in order to attenuate the contribution of  $s_i(n)$  in  $x_i(n)$  since only  $q_{ki}(n)$  for  $k \neq i$  contains information about the distance between the transceivers. Therefore,  $q_{ki}(n)$  is modelled explicitly in terms of the delay parameter (in samples)  $\eta \in [M, K]$  with  $M < K < N$ , which we are interested in estimating, and the gain  $\beta \geq 0$  as

$$q_{ki}(n) = \beta s_k(n-\eta), \quad \text{for } i \neq k. \quad (4)$$

This model describes the sound propagation of the direct path. Note that the reverberation is later modelled as part of the noise (see Sec. 2.1.2) and that  $\beta$  and  $\eta$  are not indexed since we assume that they are the same for both  $q_{12}(n)$  and  $q_{21}(n)$ .

Moreover, the delay  $\eta$  is related to the distance between the loudspeakers  $d$  via  $d = \eta c / f_s$  where  $c$  is the speed of sound and  $f_s$  the sampling frequency. If we define the vectors

$$\mathbf{x}_i = [x_i(0) \ x_i(1) \ \cdots \ x_i(N-1)]^T \quad (5)$$

$$\mathbf{x} = [\mathbf{x}_1^T \ \mathbf{x}_2^T]^T \quad (6)$$

$$\mathbf{s}_i(\eta) = [s_i(-\eta) \ s_i(1-\eta) \ \cdots \ s_i(N-1-\eta)]^T \quad (7)$$

$$\mathbf{e}_i = [e_i(0) \ e_i(1) \ \cdots \ e_i(N-1)]^T \quad (8)$$

$$\mathbf{e} = [\mathbf{e}_1^T \ \mathbf{e}_2^T]^T \quad (9)$$

$$\mathbf{h}_i = [h_i(0) \ h_i(1) \ \cdots \ h_i(M-1)]^T, \quad (10)$$

it follows that the signal model can be written as

$$\mathbf{x} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{0} & \mathbf{s}_2(\eta) \\ \mathbf{0} & \mathbf{B}_2 & \mathbf{s}_1(\eta) \end{bmatrix} \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \beta \end{bmatrix} + \mathbf{e} \quad (11)$$

$$= \mathbf{B}\mathbf{h} + \mathbf{s}(\eta)\beta + \mathbf{e} \quad (12)$$

where the definitions of  $\mathbf{B}$ ,  $\mathbf{h}$ , and  $\mathbf{s}(\eta)$  are obvious and

$$\mathbf{B}_i = [s_i(0) \ s_i(1) \ \cdots \ s_i(M-1)] \quad (13)$$

is a convolution matrix. To summarise, we have so far assumed a signal model which is linear in the unknown transceiver filters  $\mathbf{h}_1$  and  $\mathbf{h}_2$  and the gain  $\beta$  and is non-linear in the delay  $\eta$ . The main reason for using this signal model is that, as we show in Sec. 2.2, the linear parameters can easily be separated out of the problem leaving us with the single non-linear parameter  $\eta$  which we are interested in estimating. Before deriving the estimator for  $\eta$ , however, we make a number of assumptions about the source signal and the noise which enables sub-sample delay estimation accuracy, drastically reduces the computational complexity, and increases the robustness of the resulting estimator.

#### 2.1.1. The Source Signals

Most scientific literature on time of arrival (TOA), time difference of arrival (TDOA), and DOA estimation formulates these problems in the frequency domain since a delay in the time domain corresponds to a phase-shift in the frequency domain. Consequently, the delay parameter can be separated out analytically from the source signal and modelled as a continuous parameter. For finite length signals, however, a delay in the time domain only corresponds to a phase shift in the frequency domain if the signal is periodic with fundamental frequency  $2\pi/N$  radians per sample (or an integer multiple thereof) [11]. Since we here consider very long segments compared to the delay, we wish to estimate, we do not make an inappropriate error by assuming that the source signals are periodic. Thus, we have that

$$\mathbf{s}_i(\eta) = \mathbf{Z}\mathbf{A}_i\mathbf{d}(\eta) \quad (14)$$

$$\mathbf{B}_i = \mathbf{Z}\mathbf{A}_i\mathbf{F} \quad (15)$$

where we have defined

$$\mathbf{z}(\omega) = [1 \quad \exp(j\omega) \quad \cdots \quad \exp(j\omega(N-1))]^T \quad (16)$$

$$\mathbf{Z} = [\mathbf{z}(-2\pi L/N) \quad \cdots \quad \mathbf{1} \quad \cdots \quad \mathbf{z}(2\pi L/N)] \quad (17)$$

$$\mathbf{d}(\eta) = [\exp(j2\pi\eta L/N) \quad \cdots \quad 1 \quad \cdots \quad \exp(-j2\pi\eta L/N)]^T \quad (18)$$

$$\mathbf{A}_i = N^{-1} \text{diag}(\mathbf{Z}^H \mathbf{s}_i(0)) \quad (19)$$

$$\mathbf{F} = [\mathbf{d}(0) \quad \mathbf{d}(1) \quad \cdots \quad \mathbf{d}(M-1)] \quad (20)$$

Note that the time indices are symmetric around zero from  $-L$  to  $L$  where  $L = \lfloor N/2 \rfloor$ . This is necessary to ensure that the decomposition of  $\mathbf{s}_i(\eta)$  is real-valued for non-integer values of  $\eta$  [11].

### 2.1.2. The Noise

We assume that the noise is Gaussian and consists of two parts

$$\mathbf{e}_i = \mathbf{w}_i + \mathbf{v}_i \quad (21)$$

where the first part is due to reverberation and the second part is measurement noise. These two are assumed to be independent, and the measurement noise is modelled as white Gaussian noise with variance  $\sigma^2$ . We model  $\mathbf{w}_i$  as a delayed and weighted sum of the two source signals so that

$$\mathbf{w}_i = \sum_{m=2}^M (\mathbf{s}_1(\eta_{1i,m})\beta_{1i,m} + \mathbf{s}_2(\eta_{2i,m})\beta_{2i,m}) \quad (22)$$

where  $\eta_{1i,m}$  and  $\beta_{1i,m}$  are the  $m$ 'th reflection and gain from transceiver 1 to transceiver  $i$ . The summation index is running from  $m = 2$  to indicate that the first component is already included in the model via (4). We now make the critical assumption that all reflections are uncorrelated so that

$$E[\mathbf{w}_i \mathbf{w}_i^H] \approx \mathbf{0} \quad (23)$$

$$E[\mathbf{w}_i \mathbf{w}_i^H] \approx \sum_{m=2}^M E \left[ \mathbf{s}_1(\eta_{1i,m})\beta_{1i,m}^2 \mathbf{s}_1^H(\eta_{1i,m}) + \mathbf{s}_2(\eta_{2i,m})\beta_{2i,m}^2 \mathbf{s}_2^H(\eta_{2i,m}) \right] \quad (24)$$

$$\approx \gamma \sigma^2 \mathbf{Z} (\mathbf{A}_1 \mathbf{A}_1^H + \mathbf{A}_2 \mathbf{A}_2^H) \mathbf{Z}^H \quad (25)$$

where  $\gamma$  is an uninteresting scale parameter and the last expression follows from the decomposition in (14) and from

$$E \left[ \sum_{m=2}^M \mathbf{d}(\eta_{i,m})\beta_{i,m}^2 \mathbf{d}^H(\eta_{i,m}) \right] \approx \gamma \sigma^2 \mathbf{I}_N \quad (26)$$

These assumptions are hard to justify theoretically, but have been demonstrated to work well in practice [12, 13]. Under these assumptions, the covariance matrix of the noise can be written as

$$\mathbf{C} = E[\mathbf{e}\mathbf{e}^H] \approx \begin{bmatrix} \mathbf{C}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_2 \end{bmatrix} \quad (27)$$

$$\mathbf{C}_i \approx \gamma \sigma^2 \left[ \mathbf{Z} (\mathbf{A}_1 \mathbf{A}_1^H + \mathbf{A}_2 \mathbf{A}_2^H) \mathbf{Z}^H + \gamma^{-1} \mathbf{I}_N \right] \quad (28)$$

Applying the matrix inversion lemma to  $\mathbf{C}_i^{-1}$ , we obtain that

$$\mathbf{C}_i^{-1} = \sigma^{-2} \left[ \mathbf{I}_N - N^{-1} \mathbf{Z} \mathbf{Z}^H + (N^2 \gamma)^{-1} \mathbf{Z} \mathbf{Q} \mathbf{Z}^H \right] \quad (29)$$

where we have defined

$$\mathbf{Q} = \left( \mathbf{A}_1 \mathbf{A}_1^H + \mathbf{A}_2 \mathbf{A}_2^H + (N\gamma)^{-1} \mathbf{I}_N \right)^{-1} \quad (30)$$

With these, we obtain

$$\mathbf{Z}^H \mathbf{C}_i^{-1} = (\sigma^2 N \gamma)^{-1} \mathbf{Q} \mathbf{Z}^H \quad (31)$$

$$\mathbf{Z}^H \mathbf{C}_i^{-1} \mathbf{Z} = (\sigma^2 \gamma)^{-1} \mathbf{Q} \quad (32)$$

which proves to be useful later.

## 2.2. A Maximum Likelihood Estimator

The log-likelihood function pertaining to the model in (12) is given by

$$l(\mathbf{h}_1, \mathbf{h}_2, \beta, \eta, \sigma^2, \gamma) = -\frac{1}{2} \left[ \ln |\mathbf{C}| + (\mathbf{x} - \mathbf{B}\mathbf{h} - \mathbf{s}(\eta)\beta)^H \mathbf{C}^{-1} (\mathbf{x} - \mathbf{B}\mathbf{h} - \mathbf{s}(\eta)\beta) \right] \quad (33)$$

where all terms which do not depend on the unknown parameters have been ignored. Whereas the linear parameters  $\mathbf{h}$  and  $\beta$  and the noise variance  $\sigma^2$  can be separated out of the likelihood function, the scale factor  $\gamma$  cannot. Since  $\gamma$  is only a nuisance parameter, we assume that it is known and large. That is, we assume that the reverberation energy is much larger than that of the measurement noise. We have found that this works very well in practice. As seen from (30), this means that  $(N\gamma)^{-1}$  acts as a regularisation parameter.

To derive the maximum likelihood (ML) estimator for the delay  $\eta$ , we perform the following steps. Given  $\eta$  and  $\beta$ , the ML-estimate of the transceiver filters is given by

$$\hat{\mathbf{h}} = \left( \mathbf{B}^H \mathbf{C}^{-1} \mathbf{B} \right)^{-1} \mathbf{B}^H \mathbf{C}^{-1} (\mathbf{x} - \mathbf{s}(\eta)\beta), \quad (34)$$

Inserting this estimate back into the log-likelihood function in (33) and only keeping the terms which depend on  $\eta$  and  $\beta$  give the optimisation problem<sup>1</sup>

$$\hat{\beta}, \hat{\eta} = \underset{\beta \geq 0, \eta \in [M, K]}{\text{argmin}} (\mathbf{x} - \mathbf{s}(\eta)\beta)^H \mathbf{C}^{-1} \mathbf{R} (\mathbf{x} - \mathbf{s}(\eta)\beta) \quad (35)$$

where  $\mathbf{R} = \text{diag}(\mathbf{R}_1, \mathbf{R}_2)$  is a block diagonal matrix with  $\mathbf{R}_i = \mathbf{I}_N - \mathbf{B}_i (\mathbf{B}_i^H \mathbf{C}_i^{-1} \mathbf{B}_i)^{-1} \mathbf{B}_i^H \mathbf{C}_i^{-1}$ . Despite the non-negative constraint on the gain  $\beta$ , it can still be separated out of the optimisation problem by solving a KKT system of equations. The final 1D optimisation problem for the delay is

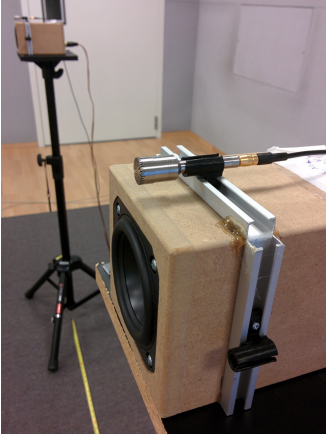
$$\hat{\eta} = \underset{\eta \in [M, K]}{\text{argmax}} \max(J(\eta), 0) \quad (36)$$

where the cost function is given by

$$J(\eta) = \frac{\mathbf{s}_2^H(\eta) \mathbf{C}_1^{-1} \mathbf{R}_1 \mathbf{x}_1 + \mathbf{s}_1^H(\eta) \mathbf{C}_2^{-1} \mathbf{R}_2 \mathbf{x}_2}{\sqrt{\mathbf{s}_2^H(\eta) \mathbf{C}_1^{-1} \mathbf{R}_1 \mathbf{s}_2(\eta) + \mathbf{s}_1^H(\eta) \mathbf{C}_2^{-1} \mathbf{R}_2 \mathbf{s}_1(\eta)}} \quad (37)$$

This cost function is highly non-linear in  $\eta$  so we propose to find  $\hat{\eta}$  using a two step procedure. First, a coarse value for  $\hat{\eta}$  is computed from a search over  $J(\eta)$  on a uniform grid. Second, the coarse estimate is refined using a line searching method such as a Fibonacci search [14, pp. 85–92].

<sup>1</sup>Note that  $\mathbf{R}^H \mathbf{C}^{-1} \mathbf{R} = \mathbf{C}^{-1} \mathbf{R}$ .



**Fig. 1.** A picture of the stereo setup.

Type	WGN	Music	Speech
Simulation	0.28	0.78	0.18
Measurement	0.61	1.42	1.13

**Table 1.** Standard deviation in mm of the estimated distance for three source signals in a simulated and real environment.

### 2.2.1. Efficient Implementation

The cost function  $J(\eta)$  can be evaluated efficiently by using the intermediate results in (14), (15), (31), and (32), and by computing the economy size singular value decomposition (SVD)  $\mathbf{Q}^{1/2}\mathbf{A}_i\mathbf{F} = \mathbf{U}_i\mathbf{S}_i\mathbf{V}_i$  so that

$$\mathbf{Z}^H\mathbf{C}_i^{-1}\mathbf{R}_i = (\sigma^2N\gamma)^{-1}\mathbf{Q}^{1/2}(\mathbf{I}_N - \mathbf{U}_i\mathbf{U}_i^H)\mathbf{Q}^{1/2}\mathbf{Z}^H.$$

These results allow us to write the cost function as

$$J(\eta) = \frac{\mathbf{d}^H(\eta)(\mathbf{y}_1 + \mathbf{y}_2)}{\sqrt{2L + 1 - \mathbf{d}^H(\eta)(\mathbf{K}_1 + \mathbf{K}_2)\mathbf{d}(\eta)}} \quad (37)$$

where (for  $k \neq i$ )

$$\mathbf{y}_i = \mathbf{A}_k^H\mathbf{Q}^{1/2}(\mathbf{I}_N - \mathbf{U}_i\mathbf{U}_i^H)\mathbf{Q}^{1/2}\mathbf{Z}^H\mathbf{x}_i \quad (38)$$

$$\mathbf{K}_i = \mathbf{A}_i^H\mathbf{Q}^{1/2}\mathbf{U}_i\mathbf{U}_i^H\mathbf{Q}^{1/2}\mathbf{A}_i. \quad (39)$$

Note that  $\mathbf{Z}^H\mathbf{x}_i$  and all elements of the diagonal matrices  $\mathbf{A}_i$  and  $\mathbf{Q}$  can be computed using an FFT algorithm. Moreover,  $\mathbf{d}^H(\eta)\mathbf{K}_i\mathbf{d}(\eta)$  is approximately zero and depends only weakly on  $\eta$  since  $\mathbf{d}(\eta)$  is asymptotically orthogonal to the columns of  $\mathbf{F}$  for  $\eta \geq M$ . Therefore, we have in practice found that only the numerator in the cost function is sufficient to find the coarse estimate of  $\eta$ . On the Fourier grid, the numerator can be computed using a single FFT whereas the denominator requires  $2M$  FFTs.

## 3. RESULTS

In this section, we demonstrate the applicability of the proposed method in both a simulated and a real environment. The former is necessary to be able to compare the produced

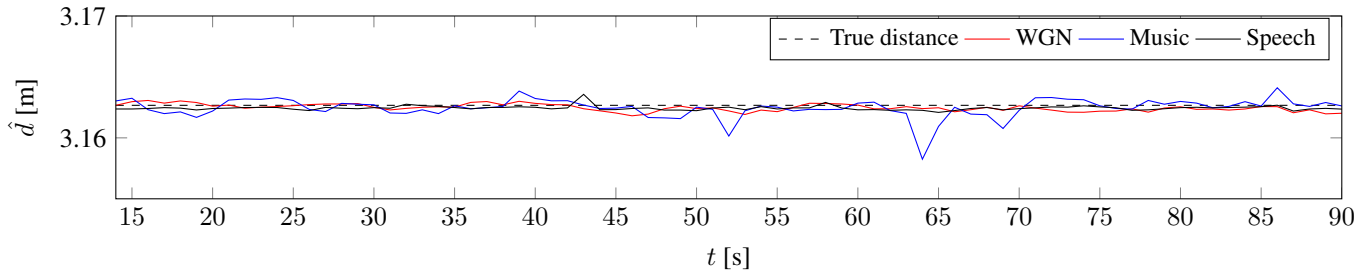
estimates to a ground truth which is unknown and not well defined in a real environment. Specifically, we evaluated the estimator for three different source signals: (1) a white Gaussian noise signal, (2) a stereo music signal (track 61 on the EBU SQAM cd [15]), and (3) a stereo speech signal (track 4 on the Archimedes CD [16]). All signals were played back and recorded at a sampling rate of 44.1 kHz. The source signals to the loudspeakers were also recorded to remove internal delays in the PC and the sound card. Data frames of four seconds were obtained with a 75 % of overlap between the successive frames. These data were down-sampled by a factor of four to 11025 Hz since the 3" loudspeakers used in the measurements and shown in Fig. 1 have a very non-linear response at the higher frequencies. A MATLAB implementation of the proposed algorithm can process this amount of data in real-time on a standard desktop PC. For this sampling frequency and a speed of sound of 343 m/s, the sampling grid corresponds to a resolution of 3.1 cm. The code for running the simulation and making the measurements is available from <http://kom.aau.dk/~jkn/publications/publications.php>.

Fig. 2 shows an excerpt of the results of the simulation where the sources were assumed to be point sources and artificial reverberation [17] was added with a reverberation time of 0.5 seconds. From the figure and Table 1, we see that we got sub-millimetre accuracy for all source signals. From Fig. 3 and Table 1, we see that the variation of the estimates increased in the real environment despite that the loudspeakers were closer together. The main reason for this is that loudspeakers are not omnidirectional point sources. Instead, especially the higher frequencies are attenuated from one loudspeaker to the other when the loudspeakers are configured in a stereo setup as in Fig. 1, i.e., they are not pointed towards each other. Moreover, the acoustic centre of the loudspeaker is typically in front of the loudspeaker and frequency dependent [18].

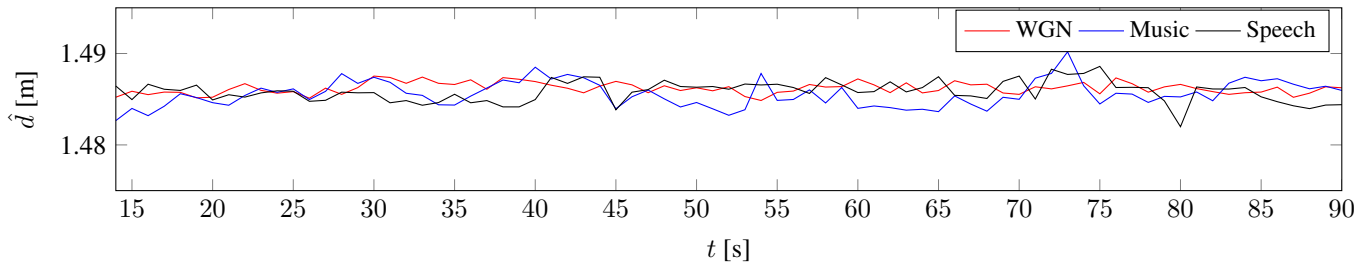
Although not shown here, outliers in the estimated distances occur occasionally. These happen typically in very silent parts of the music/speech and can be removed by using a sound activity detector or by post-processing the computed estimates using a smoothing algorithm. However, even without these heuristics, we are able to estimate the transceiver distance to a millimetre precision for even a modest sampling frequency.

## 4. CONCLUSION

We have here taken a first step in the direction of using music or speech signals for the localisation of a number of loudspeakers. Specifically, we have considered the simplest case where we have to estimate the distance between two loudspeakers, each equipped with a single microphone. We derived an ML estimator for this problem and demonstrated that it could be used to obtain real-time distance estimates to



**Fig. 2.** The estimated distance for three source signals in a simulated environment.



**Fig. 3.** The estimated distance for three source signals in a real environment.

within an accuracy of one millimetre for even a low sampling frequency. Only frame-by-frame processing was considered in this paper, but outliers can be removed and higher accuracy can be achieved by smoothing the computed estimates.

## 5. REFERENCES

- [1] F. Rumsey, D. Griesinger, T. Holman, M. Sawaguchi, G. Steinke, G. Theile, and T. Wakatuki, "Multichannel surround sound systems and operations," The Audio Engineering Society, Tech. Rep. AESTD1001.0.01-05, 2001.
- [2] C. Kyriakakis, "Fundamental and technological limitations of immersive audio systems," *Proc. IEEE*, vol. 86, no. 5, pp. 941–951, 1998.
- [3] B. Gunel, "Loudspeaker localization using B-format recordings," in *Proc. IEEE Workshop on Appl. of Signal Process. to Aud. and Acoust.* IEEE, 2003, pp. 59–62.
- [4] B. Atkinson, T. Blank, M. Isard, J. D. Johnston, and K. Olynyk, "An internet protocol (IP) sound system," in *Proc. AES Convention*. Audio Engineering Society, 2004.
- [5] S. Choisel, G. G. Martin, and M. Hlatky, "Loudspeaker position estimation," US Patent 8 279 709 B2, 2012.
- [6] G. Del Galdo, M. Lang, J. A. Pineda Pardo, A. Silzle, and O. Thiergart, "Acoustic measurement system for 3-D loudspeaker set-ups," in *Audio Eng. Soc.* Audio Engineering Society, 2010.
- [7] A. J. V. Leest and D. W. Schobben, "Method of and system for determining distances between loudspeakers," US Patent 7 864 631 B2, 2011.
- [8] F. Kolbeck, G. Del Galdo, I. Sobieraj, and T. Bliem, "Loudspeaker localization based on audio watermarking," in *Proc. AES Convention*. Audio Engineering Society, 2012.
- [9] W. S. Torgerson, "Multidimensional scaling: I. theory and method," *Psychometrika*, vol. 17, no. 4, pp. 401–419, 1952.
- [10] M. Crocco, A. Del Bue, and V. Murino, "A bilinear approach to the position self-calibration of multiple sensors," *IEEE Trans. Signal Process.*, vol. 60, no. 2, pp. 660–673, 2012.
- [11] J. R. Jensen, J. K. Nielsen, M. G. Christensen, and S. H. Jensen, "On frequency domain models for TDOA estimation," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015.
- [12] C. Zhang, D. Florencio, and Z. Zhang, "Why does PHAT work well in low noise, reverberative environments?" in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* IEEE, 2008, pp. 2565–2568.
- [13] Y. Rui and D. Florencio, "Time delay estimation in the presence of correlated noise and reverberation," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, vol. 2. IEEE, 2004, pp. 133–136.
- [14] A. Antoniou and W.-S. Lu, *Practical Optimization: Algorithms and Engineering Applications*. Springer, Mar. 2007.
- [15] E. B. Union, "Sound quality assessment material recordings for subjective tests: Users' handbook for the EBU SQAM CD," European Broadcasting Union, Tech. Rep. EBU – TECH 3253, 2008.
- [16] Bang & Olufsen. (1992) Music for archimedes. Compact disc CD B&O 101.
- [17] E. A. P. Habets. (2010) Room impulse response generator. Ver. 2.0.20100920. [Online]. Available: <http://www.audiolabs-erlangen.de/fau/professor/habets/software/rir-generator>
- [18] F. Jacobsen, S. B. Figueroa, and K. Rasmussen, "A note on the concept of acoustic center," *J. Acoust. Soc. Am.*, vol. 115, no. 4, pp. 1468–1473, 2004.