

# Rectified binaural ratio: A complex T-distributed feature for robust sound localization

Antoine Deleforge, Florence Forbes

► **To cite this version:**

Antoine Deleforge, Florence Forbes. Rectified binaural ratio: A complex T-distributed feature for robust sound localization. European Signal Processing Conference, Aug 2016, Budapest, Hungary. IEEE, Proceedings of the 24th European Signal Processing Conference (EUSIPCO), 2016, pp.1257-1261, 2016, <10.1109/EUSIPCO.2016.7760450>. <hal-01372337>

**HAL Id: hal-01372337**

**<https://hal.inria.fr/hal-01372337>**

Submitted on 27 Sep 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# RECTIFIED BINAURAL RATIO: A COMPLEX T-DISTRIBUTED FEATURE FOR ROBUST SOUND LOCALIZATION

Antoine Deleforge\* and Florence Forbes†

\*Inria Rennes - Bretagne Atlantique †Inria Grenoble - Rhône-Alpes (firstname.lastname@inria.fr)

## ABSTRACT

Most existing methods in binaural sound source localization rely on some kind of aggregation of phase- and level- difference cues in the time-frequency plane. While different aggregation schemes exist, they are often heuristic and suffer in adverse noise conditions. In this paper, we introduce the *rectified binaural ratio* as a new feature for sound source localization. We show that for Gaussian-process point source signals corrupted by stationary Gaussian noise, this ratio follows a complex t-distribution with explicit parameters. This new formulation provides a principled and statistically sound way to aggregate binaural features in the presence of noise. We subsequently derive two simple and efficient methods for robust relative transfer function and time-delay estimation. Experiments on heavily corrupted simulated and speech signals demonstrate the robustness of the proposed scheme.

**Index Terms**— Complex Gaussian ratio; t-distribution; relative transfer function; binaural; sound localization

## 1. INTRODUCTION

The most widely used features for binaural (two microphones) sound source localization are the measured time delays and level differences between the two microphones. For a single source signal in the absence of noise, these features correspond in the frequency domain to the ratio of the Fourier transforms of the right- and the left-microphone signals. This ratio is called the *relative transfer function* (RTF) [1], and only depends on the source’s spatial characteristics, *e.g.*, its position relative to the microphones. The log-amplitudes and phases of the RTF are referred to as *interaural level differences* (ILD) and *interaural phase differences* (IPD) in the binaural literature. Many binaural sound source localization methods rely on some kind of aggregation of these cues over the time-frequency plane [2–8]. The generalized cross-correlation (GCC) method [2] consists of weighting the cross-power spectral density (CPSD) of two signals in order to estimate their delay in the time-domain (CPSD phases and IPD are the same). A successful GCC method is the phase transform (PHAT), in which IPD cues are equally weighted. The popular sound localization method PHAT-histogram aggregates these cues using histograms [3]. In [5], a heuristic binaural cue weighting scheme based on

signals’ onsets is proposed. In [4], both ILD and IPD cues are modeled as real Gaussians and their frequency-dependent variances are estimated through an expectation-maximization (EM) procedure referred to as MESSL. A number of extensions of MESSL have later been developed [6–8], including one using t-distributions for ILD and IPD cues instead of Gaussian distributions [6].

While all these methods rely on a weighting scheme of binaural cues, none of these schemes is based on the statistical properties of the source and noise signals. Though, intuitively, a low signal-to-noise-ratio (SNR) at microphones means that a specific cue is less reliable, while a high SNR means that this cue should be given more weight. In this paper, we prove that the ratio of two complex circular-symmetric Gaussian variables follows a complex t-distribution with explicit parameter expressions. In particular, for the binaural recording of a Gaussian-process source corrupted by stationary Gaussian noise, we show that the mean of the microphone signals’ ratio does not only depend on the clean ratio but also on the source and noise statistics. This observation naturally leads to the definition of a new binaural feature referred to as the *rectified binaural ratio* (RBR). The explicit distribution of RBR features provides a principled and statistically sound way of weighting and aggregating them. Based on this, we derive two simple and efficient methods for relative transfer function and time-delay estimation, and test their robustness on heavily corrupted binaural signals.

## 2. A COMPLEX-T MODEL FOR BINAURAL CUES

In the complex short-time Fourier domain, we consider the following model for a binaural setup recording a static point sound source in the presence of noise:

$$\begin{cases} m_1(f, t) = h_1(f, \boldsymbol{\theta})s(f, t) + n_1(f, t) \\ m_2(f, t) = h_2(f, \boldsymbol{\theta})s(f, t) + n_2(f, t) \end{cases},$$

or equivalently  $\mathbf{m}(f, t) = \mathbf{h}(f, \boldsymbol{\theta})s(f, t) + \mathbf{n}(f, t)$ . (1)

Here,  $(f, t)$  is the frequency-time indexing,  $\boldsymbol{\theta}$  is a vector of source spatial parameters, *e.g.*, the source position,  $\mathbf{m}(f, t) = [m_1(f, t), m_2(f, t)]^T \in \mathbb{C}^2$  denotes the microphone signals,  $s(f, t) \in \mathbb{C}$  denotes the source signal of interest,  $\mathbf{n}(f, t) = [n_1(f, t), n_2(f, t)]^T \in \mathbb{C}^2$  denotes the noise signals and  $\mathbf{h}(f, \boldsymbol{\theta}) = [h_1(f, \boldsymbol{\theta}), h_2(f, \boldsymbol{\theta})]^T \in \mathbb{C}^2$  denotes

the acoustic transfer function from the source to the microphones. The function  $\mathbf{h}(f, \boldsymbol{\theta})$  is of particular interest because it depends on the source position  $\boldsymbol{\theta}$  but does not depend on the time-varying source and noise signals. Under noise-free and non-vanishing source assumptions, *i.e.*  $\mathbf{n}(f, t) = \mathbf{0}$  and  $s(f, t) \neq 0$ , it is easily seen that the *binaural ratio*  $m_2(f, t)/m_1(f, t)$  is equal to  $h_2(f, \boldsymbol{\theta})/h_1(f, \boldsymbol{\theta}) = r(f, \boldsymbol{\theta})$ , which only depends on the source position. This ratio can hence be used for sound source localization. The quantity  $r(f, \boldsymbol{\theta})$  is called *relative transfer function* (RTF) [1]. Its log-amplitudes and phases are respectively referred to as interaural level and phase differences (ILD and IPD).

In practical situations including noise, the ratio  $m_2(f, t)/m_1(f, t)$  does no longer depend on  $\boldsymbol{\theta}$  only, but also on the source and noise signals  $s(f, t)$  and  $\mathbf{n}(f, t)$ . These signals are assumed independent, and we consider the following probabilistic models:

$$\begin{aligned} P(s(f, t)) &= \mathcal{CN}_1(s(f, t); 0, \sigma_s^2(f, t)), \\ P(\mathbf{n}(f, t)) &= \mathcal{CN}_2(\mathbf{n}(f, t); \mathbf{0}, \mathbf{R}_{nn}(f)), \end{aligned} \quad (2)$$

where  $\mathcal{CN}_p$  denotes the  $p$ -variate complex circular-symmetric normal distribution, or *complex-normal*. Its density is [9]:

$$\mathcal{CN}_p(\mathbf{x}; \mathbf{c}, \boldsymbol{\Sigma}) = \frac{1}{\pi^p |\boldsymbol{\Sigma}|} \exp(-(\mathbf{x} - \mathbf{c})^H \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \mathbf{c})),$$

where  $\{\cdot\}^H$  denotes the Hermitian transpose. We assume that  $\mathbf{R}_{nn}(f)$  is known and constant over time, *i.e.*, noise signals are stationary. However, they are not necessarily pairwise independent and may thus include other point sources. On the other hand, the source signal is a Gaussian process with time-varying variance  $\sigma_s^2(f, t)$ . This general model is widely used in audio signal processing, in particular for sound source separation, *e.g.*, [10]. We now introduce the univariate *complex t-distribution* denoted  $\mathcal{CT}_1$ :

$$\mathcal{CT}_1(y; \mu, \lambda^2, \nu) = \frac{1}{\pi \lambda^2} \left( 1 + \frac{|y - \mu|^2}{\nu \lambda^2} \right)^{-(1+\nu)}, \quad (4)$$

where  $\mu \in \mathbb{C}$ ,  $\lambda^2 \in \mathbb{R}^+$  and  $\nu \in \mathbb{R}^+$  are respectively referred to as the mean, spread and degrees of freedom parameters. This definition follows a construction of multivariate extensions for the t-distribution [11] applied to the complex plane. In the real case, the t-distribution arises from the ratio of a Gaussian over the square root of a Chi-square distribution. In the complex case, we alternatively show the following result:

**Theorem 1** *Let  $\mathbf{m} = [m_1, m_2]^\top$  be a vector in  $\mathbb{C}^2$  following a complex-normal distribution such that*

$$P(\mathbf{m}) = \mathcal{CN}_2\left(\mathbf{m}; \mathbf{0}, \begin{bmatrix} \sigma_{m_1}^2 & \rho \sigma_{m_1} \sigma_{m_2} \\ \rho^* \sigma_{m_1} \sigma_{m_2} & \sigma_{m_2}^2 \end{bmatrix}\right).$$

*Then the ratio variable  $y = m_2/m_1$  follows a complex-t distribution such that*

$$P(y) = \mathcal{CT}_1\left(y; \frac{\sigma_{m_2}}{\sigma_{m_1}} \rho^*, \frac{\sigma_{m_2}^2}{\sigma_{m_1}^2} (1 - |\rho|^2), 1\right). \quad (5)$$

Here,  $\rho = \mathbb{E}\{m_1 m_2^*\} / (\sigma_{m_1} \sigma_{m_2})$  is the correlation coefficient between  $m_1$  and  $m_2$  and  $(\cdot)^*$  denotes the complex conjugate. This result is consistent with that in [12] but we provide a simpler proof with better insight in Appendix A.2. Theorem 1 can be directly applied to obtain an explicit distribution for the binaural ratio  $m_2(f, t)/m_1(f, t)$  under the model defined by (1), (2) and (3). However, both the mean and the spread of this distribution depend on the noise correlation and variances as well as the transfer functions in a way which is difficult to handle. We will therefore design a more convenient and somewhat more natural binaural feature by first *whitening* the noise signals in each observed vectors  $\mathbf{m}(f, t)$ , *i.e.*, making them independent and of unit variance. Since  $\mathbf{R}_{nn}(f)$  is positive semi-definite, it has a unique positive semi-definite square root  $\mathbf{R}_{nn}(f)^{1/2}$ . If  $\mathbf{R}_{nn}(f)$  is further invertible<sup>1</sup>, we can define:

$$\mathbf{Q}(f) = \mathbf{R}_{nn}(f)^{-1/2}. \quad (6)$$

By left-multiplication of (1) by  $\mathbf{Q}(f)$  we obtain

$$\mathbf{Q}(f) \mathbf{m}(f, t) = \mathbf{Q}(f) \mathbf{h}(f, \boldsymbol{\theta}) s(f, t) + \mathbf{Q}(f) \mathbf{n}(f, t), \quad (7)$$

$$\mathbf{m}'(f, t) = \mathbf{h}'(f, \boldsymbol{\theta}) s(f, t) + \mathbf{n}'(f, t), \quad (8)$$

where  $\mathbf{n}'(f, t)$  follows the standard bivariate complex-normal  $\mathcal{CN}_2(\mathbf{0}, \mathbf{I}_2)$ . Note that  $\mathbf{h}'(f, \boldsymbol{\theta})$  can only be identified up to a multiplicative complex scalar constant because the same observations are obtained by dividing corresponding source signals by this constant. Hence, we can assume without loss of generality that  $h'_1(f, \boldsymbol{\theta}) = 1$  and  $h'_2(f, \boldsymbol{\theta}) = r'(f, \boldsymbol{\theta})$ , where  $r'(f, \boldsymbol{\theta})$  is the relative transfer function (RTF) after whitening. It follows that,  $m'_1(f, t) = s(f, t) + n'_1(f, t)$ ,  $\sigma_{m'_1}^2(f, t) = \sigma_s^2(f, t) + 1$  and  $\sigma_{m'_2}^2(f, t) = |r'|^2 \sigma_s^2(f, t) + 1$ . Moreover, since  $\mathbf{Q}(f)$  is invertible, the original RTF can be obtained from  $r'(f, \boldsymbol{\theta})$  as the ratio of vector  $\mathbf{Q}(f)^{-1} [1, r'(f, \boldsymbol{\theta})]^\top$ .

We can now use Theorem 1 to obtain that  $y'(f, t) = m'_2(f, t)/m'_1(f, t)$  follows the complex-t distribution:

$$\mathcal{CT}_1\left(\frac{\sigma_s^2(f, t)}{1 + \sigma_s^2(f, t)} r'(f, \boldsymbol{\theta}), \frac{\sigma_{m'_2}^2(f, t) + \sigma_s^2(f, t)}{(1 + \sigma_s^2(f, t))^2}, 1\right). \quad (9)$$

Interestingly, it turns out that the distribution of a binaural ratio under white Gaussian noise is *not* centered on the actual RTF  $r'(f, \boldsymbol{\theta})$ ; but rather on a scaled version of it which depends on the instantaneous source variance. This suggests to use the following more natural feature that we refer to as *rectified binaural ratio* (RBR):

$$y(f, t) = \frac{1 + \sigma_s^2(f, t)}{\sigma_s^2(f, t)} \cdot \frac{m'_2(f, t)}{m'_1(f, t)}. \quad (10)$$

This feature has the following distribution:

$$P(y(f, t)) = \mathcal{CT}_1(y(f, t); r'(f, \boldsymbol{\theta}), \lambda^2(f, t), 1), \quad (11)$$

$$\text{where } \lambda^2(f, t) = \frac{\sigma_{m'_2}^2(f, t) + \sigma_s^2(f, t)}{\sigma_s^4(f, t)}, \quad (12)$$

<sup>1</sup>For the case where  $\mathbf{R}_{nn}(f)$  is non-invertible, see Appendix A.1.

which is centered on the RTF  $r'(f, \theta)$ . The spread parameter  $\lambda^2(f, t)$  is also important because it models the uncertainty or “reliability” associated to each RBR feature: the larger is  $\lambda^2(f, t)$ , the less reliable is  $y(f, t)$ . Since the noise variance is fixed to 1, we see in (12) that  $\lambda^2(f, t)$  tends to 0 when the SNR at  $(f, t)$  tends to infinity, while  $\lambda^2(f, t)$  tends to infinity when the SNR approaches 0, which matches intuition.

### 3. PARAMETER ESTIMATION

#### 3.1. Spread parameter

We consider the general case of time-varying source variances  $\sigma_s^2(f, t)$ . This is more challenging than a stationary model but also more realistic since typical audio signals such as speech or music are often sparse and impulsive in the time-frequency plane. In this case, the calculation of RBR features (10) and of their spread parameter (12) requires the knowledge of instantaneous source and microphone variances at each  $(f, t)$ . A number of ways can be envisioned to estimate them. In this paper, we use the perhaps most straightforward approach: the instantaneous microphone variances  $\sigma_{m_1}^2(f, t)$  and  $\sigma_{m_2}^2(f, t)$  are approximated by their observed magnitudes  $|m_1'(f, t)|^2$  and  $|m_2'(f, t)|^2$ . More accurate estimates could be obtained using, *e.g.*, a sliding averaging window in the time-frequency plane as in [10]. However, this simple scheme showed good performance in practice. It leads to the following straightforward estimate for  $\sigma_s^2(f, t)$ :

$$\hat{\sigma}_s^2(f, t) = \begin{cases} |m_1'(f, t)|^2 - 1 & \text{if } |m_1'(f, t)|^2 > 1, \\ 0 & \text{otherwise,} \end{cases} \quad (13)$$

from which we deduce  $\hat{\lambda}^2(f, t)$  using (12).  $\hat{\sigma}_s^2(f, t) = 0$  leads to  $\hat{\lambda}^2(f, t) = +\infty$ , corresponding to a missing data at  $(f, t)$ .

#### 3.2. Unconstrained RTF

Once the spread parameter is estimated, we are left with the estimation of  $r'(f, \theta)$  which is the mean of the complex t-distribution (12). The equivalent characterization of the t-distribution as a Gaussian scale mixture leads naturally to an EM algorithm that converges under mild conditions to the maximum likelihood [13]. Introducing an additional set of latent variables  $\mathbf{u} = \{u(f, t), f = 1 : F, t = 1 : T\}$ , we can write (11) equivalently as:

$$P(y(f, t)|u(f, t)) = \mathcal{CN}_1(y(f, t); r'(f, \theta), \frac{\lambda^2(f, t)}{u(f, t)}), \quad (14)$$

$$P(u(f, t)) = \mathcal{G}(1, 1), \quad (15)$$

where  $\mathcal{G}$  denotes the Gamma distribution. At each iteration  $(q)$ , the M-step updates  $r'(f, \theta)$  as a weighted sum of the  $y(f, t)$ 's while the E-step consists of updating the weights defined as  $\omega_{ft}^{(q)} = \frac{1}{2} \hat{\lambda}^{-2}(f, t) \cdot \mathbb{E}[u(f, t)|y(f, t); r^{(q)}(f, \theta)]$ :

$$\text{M-step: } r^{(q+1)}(f, \theta) = (\sum_{t=1}^T \omega_{ft}^{(q)} y(f, t)) / (\sum_{t=1}^T \omega_{ft}^{(q)}),$$

$$\text{E-step: } \omega_{ft}^{(q+1)} = \left( \hat{\lambda}^2(f, t) + |y(f, t) - r^{(q+1)}(f, \theta)|^2 \right)^{-1}.$$

The initial weights  $\omega_{ft}^{(0)}$  can be set to 1, although our experiments showed that random initializations usually converged to the same solution. Convergence is assumed reached when  $r'(f, \theta)$  varies by less than 0.1% at a given iteration. In practice, the algorithm converged in less than 100 iterations in nearly all of our experiments. Once an estimate  $\hat{r}'(f, \theta)$  is obtained, the non-whitened RTF  $\hat{r}(f, \theta)$  is calculated as the ratio of vector  $\mathbf{Q}(f)^{-1}[1, \hat{r}'(f, \theta)]^\top$ .

#### 3.3. Acoustic space prior on the RTF

In practice, when a sound source emits in a real room, the RTF can only take a restricted set of values belonging to the so-called *acoustic space manifold* of the system [8]. Hence, a common approach is to search for the optimal  $r'$  among a finite set of  $K$  possibilities corresponding to different locations of the source, namely  $r' \in \mathcal{R}' = \{r'_1, \dots, r'_K\}$  where  $r'_k(f) = r'(f, \theta_k)$ . From a Bayesian perspective, this corresponds to a mixture-of-Dirac prior on  $r'(f, \theta)$ . Considering the observed features  $y$ , we then look for the  $r'_k$  that maximizes the log-likelihood of  $y$  as induced by (11). Taking the logarithm of (4), this amounts to minimize:

$$\hat{k} = \underset{k=1:K}{\operatorname{argmin}} \sum_{t=1}^T \sum_{f=1}^F \log(\hat{\lambda}^2(f, t) + |y(f, t) - r'_k(f)|^2).$$

We recover the robustness property that a data point with high spread has less impact on the estimation of  $r'$ .

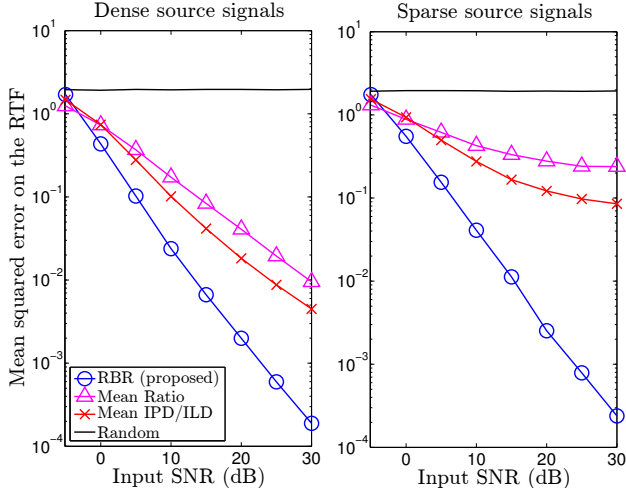
## 4. EXPERIMENTAL RESULTS

#### 4.1. RTF estimation

We first evaluate the RTF estimation method described in Section 3.2 through extensive simulations. 160,000 binaural test signals are generated according to model (1), (2) and (3), under a wide range of noise and source statistics. Each generated complex signal corresponds to  $T = 20$  time samples in a given frequency. The variances of source signals are time-varying and uniformly drawn at random. Sparse source signals are simulated by setting their variance to 0 with a 50% probability at each sample. For each test signal, the noise variances and correlation are uniformly drawn at random, and the RTF  $r$  is drawn from a standard complex-normal distribution. The proposed method is compared to two baseline methods. The first one (Mean ratio) takes the mean of the complex microphone ratios  $m_2(f, t)/m_1(f, t)$  over the  $T$  samples of each signal. The second one (Mean ILD/IPD) calculates the mean ILD and IPD as follows:

$$\begin{cases} \overline{\text{ILD}} = \frac{1}{T} \sum_{t=1}^T \log \left( \frac{m_2(f, t)}{m_1(f, t)} \right), \\ \overline{\text{IPD}} = \frac{1}{T} \sum_{t=1}^T \frac{m_2(f, t)/|m_2(f, t)|}{m_1(f, t)/|m_1(f, t)|}. \end{cases} \quad (16)$$

The RTF is then estimated as  $\exp(\overline{\text{ILD}}) \cdot \overline{\text{IPD}}$ . This latter type of binaural cue aggregation is common to many methods, including [3, 5, 8]. For fairness of comparison, the samples identified as *missing* by our method according to (13) are ignored

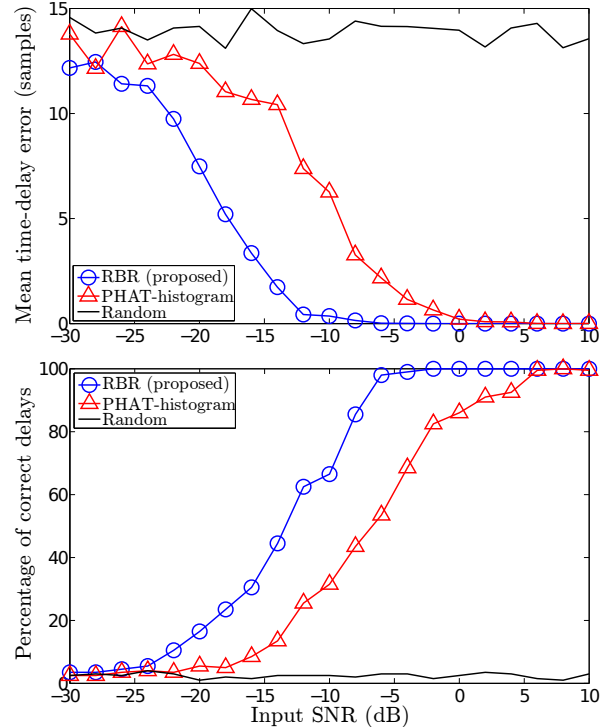


**Fig. 1.** Mean squared error of different RTF estimation methods for various SNRs.

by all 3 methods. Mean squared errors for various signal-to-noise ratios (SNR) and for both dense (left) and sparse (right) source signals are showed in Fig. 1. As an indicator of the error upper-bound, the results of a method generating random RTF estimates (Random) are also shown. Except at low SNRs ( $\leq -5$ dB) where all 3 methods yield estimates close to randomness, the proposed method outperforms both the others. In particular, for SNRs larger than 15 dB, the mean squared error is decreased by several orders of magnitudes and the RBR features performed best in 92% of the tests. Two facts may explain these results. First, as showed in (9), the microphone ratio is a biased estimate of the RTF under white noise conditions. This bias is further amplified for arbitrary noise statistics. Second, the baseline methods, as many existing methods in the literature, aggregate binaural cues with binary weights: each sample is classified as either missing or not. In contrast, the explicit spread parameter (12) available for rectified binaural ratios enables to weight observations in a statistically sound way.

#### 4.2. Time difference of arrival estimation

Under free-field conditions, *i.e.*, direct single-path propagation from the sound source to the microphones, localizing the source is equivalent to estimating the time difference of arrival (TDOA) between microphones. Indeed, for far enough sources, we have the relation  $\tau \approx d \cos(\theta) F_s / C$  where  $\tau$  is the delay in samples,  $d$  the inter-microphone distance,  $\theta$  the source's azimuth angle,  $F_s$  the frequency of sampling, and  $C$  the speed of sound. In the frequency domain, the RTF then has the explicit expression  $r(f, \tau) = \exp(-2\pi i \tau (f - 1) / F)$  where  $F$  is the number of positive frequencies and  $f = 1 : F$  is the frequency index. Let  $\mathcal{R}$  be the discrete set of RTFs corresponding to delays of  $-\tau_{\max}$  to  $+\tau_{\max}$  samples, and  $\mathcal{R}'$  the corresponding set after whitening, *i.e.*, containing ratios of  $\mathbf{Q}(f)[1, r(f, \tau)]^T$ . Given a noisy binaural signal, the method of Section 3.3 can be applied to select the most likely RTF  $r'$



**Fig. 2.** Comparing time-delay estimation results of RBR and PHAT using 1 second noisy speech signals (200 test signals per SNR value).

in  $\mathcal{R}'$  and deduce the corresponding TDOA. 4,000 test signals are generated using random 1 second speech utterances from the TIMIT dataset [14] sampled at  $F_s = 16,000$  Hz. A binaural signal with a random delay of  $-20$  to  $+20$  samples between microphones is generated, before applying the short-time Fourier transform (64ms windows with 50% overlap). This yields  $F = 512$  positive frequencies and  $T = 32$  time samples. These signals are finally corrupted by random additive stationary noise of known statistics in the frequency domain using the same procedure as in Section 4.1. The proposed RBR-based approach is compared to the sound source localization method PHAT-histogram<sup>2</sup> [3]. Results are displayed in Fig. 2. For SNRs higher than  $-6$  dB, the proposed RBR method yields less than 0.4% incorrect delays, versus 10.1% for PHAT-histogram on the same signals. RBR's average computational time is  $80 \pm 6$ ms per second of signal on a common laptop, which is about 3 times faster than PHAT-histogram using our Matlab implementations.

## 5. CONCLUSION

We explicitly expressed the probability density function of the ratio of two microphone signals in the frequency domain in the presence of a Gaussian-process point source corrupted by stationary Gaussian noise. This statistical framework enabled us to model the uncertainty of binaural cues and was efficiently applied to robust RTF and TDOA estimation. Future

<sup>2</sup>We used the PHAT-histogram implementation of Michael Mandel, available at <http://blog.mr-pc.org/2011/09/14/messl-code-online/>.

work will include extensions to multiple sound source separation and localization following ideas in [4], and to more than two microphones following ideas in [15]. The flexibility of the proposed framework may also allow the inclusion of a variety of priors on the RTFs such as Gaussian mixtures, as well as the handling of various types of noise and source statistics.

## A. APPENDIX

### A.1. Non-invertible noise covariance

If the noise signals  $n_1(f, t)$  and  $n_2(f, t)$  in (1) have a deterministic dependency,  $\mathbf{R}_{nn}(f)$  is rank-1 and non-invertible. This is an important special case which may occur in practice when, *e.g.*, the noise is a point source. Since  $\mathbf{R}_{nn}(f)^{-1/2}$  is then not defined, we replace the *whitening* matrix in (6) by  $\mathbf{Q}(f) = \begin{bmatrix} 1/\sigma_{n_1}^2(f) & 0 \\ 1/\sigma_{n_2}^2(f) & -1/\sigma_{n_2}^2(f) \end{bmatrix}$ , where  $\sigma_{n_1}^2(f)$  and  $\sigma_{n_2}^2(f)$  denote the variances of  $n_1(f, t)$  and  $n_2(f, t)$ . It then follows that  $n'_2(f, t) = 0$  and that  $n'_1(f, t)$  follows the standard complex-normal distribution  $\mathcal{CN}(0, 1)$ . All subsequent derivations in the paper remain unchanged, with the exception of (12) which becomes  $\lambda^2(f, t) = \sigma_{m'_2}^2/\sigma_s^4$ .

### A.2. Proof of Theorem 1

We first prove the result for  $\rho = 0$ , *i.e.*, when  $m_1$  and  $m_2$  are independent. Since  $[m_1, m_2]^T$  is jointly circular symmetric complex Gaussian, it follows that  $m_1$  and  $m_2$  are also complex Gaussian with  $m_1 \sim \mathcal{CN}_1(0, \sigma_{m_1}^2)$ ,  $m_2 \sim \mathcal{CN}_1(0, \sigma_{m_2}^2)$  [16], and  $S^2 = 2|m_1|^2/\sigma_{m_1}^2$  follows a Chi-square distribution with 2 degrees of freedom [9]. These properties generalize their counterparts in the real case and can be easily checked by using the characterization of complex Gaussians as real Gaussians on the real and imaginary parts [16]. We can now use the property of circular symmetric Gaussians that states that if  $Y$  is  $\mathcal{CN}(0, \Sigma)$  then  $Y$  and  $Ye^{i\phi}$  have the same distribution for all  $\phi$ . We deduce from this property that  $y = m_2/m_1$  and  $z = m_2/|m_1|$  have the same distribution. Then,  $\sigma_{m_1}z = m_2\sqrt{2/S^2}$  is distributed as a complex Gaussian over the square root of an independent scaled Chi-square distribution, which is one of the characterization of the complex t-distribution [11, Section 5.12]. It follows that  $\sigma_{m_1}z \sim \mathcal{CT}_1(0, \sigma_{m_2}^2, 1)$ . Therefore  $y$  follows  $\mathcal{CT}_1(0, \sigma_{m_2}^2/\sigma_{m_1}^2, 1)$  which corresponds to Theorem 1 for  $\rho = 0$ . For the general case, we multiply  $\mathbf{m}$  by matrix  $\mathbf{A} = \begin{bmatrix} 1 & 0 \\ -\rho^* & \sigma_{m_1}/\sigma_{m_2} \end{bmatrix}$  so that  $\tilde{\mathbf{m}} = \mathbf{A}\mathbf{m}$  is complex Gaussian with covariance matrix  $\mathbf{A}\Sigma\mathbf{A}^H = \sigma_{m_1}^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 - |\rho|^2 \end{bmatrix}$ . We deduce from the previous case that  $\tilde{y} = \tilde{m}_2/\tilde{m}_1$  follows  $\mathcal{CT}_1(0, 1 - |\rho|^2, 1)$ . We finally obtain Theorem 1's result by noting that  $\tilde{y} = (\sigma_{m_1}/\sigma_{m_2})y - \rho^*$ .

## REFERENCES

[1] S. Gannot, D. Burshtein, and E. Weinstein, "Signal enhancement using beamforming and nonstationarity with applica-

tions to speech," *IEEE Trans. Signal Process.*, vol. 49, no. 8, pp. 1614–1626, 2001.

- [2] C. Knapp and G. C. Carter, "The generalized correlation method for estimation of time delay," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 24, no. 4, pp. 320–327, 1976.
- [3] P. Aarabi, "Self-localizing dynamic microphone arrays," *IEEE Trans. Syst., Man, Cybern. C*, vol. 32, no. 4, pp. 474–484, 2002.
- [4] M. I. Mandel, R. J. Weiss, and D. P. Ellis, "Model-based expectation-maximization source separation and localization," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 18, no. 2, pp. 382–394, 2010.
- [5] J. Woodruff and D. Wang, "Binaural localization of multiple sources in reverberant and noisy environments," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 20, no. 5, pp. 1503–1512, 2012.
- [6] Z. Zohny and J. Chambers, "Modelling interaural level and phase cues with student's t-distribution for robust clustering in MESSL," in *International Conference on Digital Signal Processing (DSP)*. IEEE, 2014, pp. 59–62.
- [7] M. I. Mandel and N. Roman, "Enforcing consistency in spectral masks using markov random fields," in *EUSIPCO*. IEEE, 2015, pp. 2028–2032.
- [8] A. Deleforge, F. Forbes, and R. Horaud, "Acoustic space learning for sound-source separation and localization on binaural manifolds," *International journal of neural systems*, vol. 25, no. 01, pp. 1440003, 2015.
- [9] D. R. Fuhrmann, "Complex random variables and stochastic processes," *The Digital Signal Processing Handbook*, pp. 60–1, 1997.
- [10] E. Vincent, S. Arberet, and R. Gribonval, "Underdetermined instantaneous audio source separation via local gaussian modeling," in *Independent Component Analysis and Signal Separation*, pp. 775–782. Springer, 2009.
- [11] S. Kotz and S. Nadarajah, *Multivariate t Distributions and their Applications*, Cambridge, 2004.
- [12] R. J. Baxley, B. T. Walkenhorst, and G. Acosta-Marum, "Complex Gaussian ratio distribution with applications for error rate calculation in fading channels with imperfect CSI," in *Global Telecommunications Conference (GLOBECOM)*. IEEE, 2010, pp. 1–5.
- [13] G. McLachlan and D. Peel, "Robust mixture modelling using the T distribution," *Statistics and computing*, vol. 10, pp. 339–348, 2000.
- [14] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, and D. S. Pallett, "The DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM," Tech. Rep. NISTIR 4930, National Institute of Standards and Technology, Gaithersburg, MD, 1993.
- [15] A. Deleforge, S. Gannot, and W. Kellermann, "Towards a generalization of relative transfer functions to more than one source," in *EUSIPCO*. IEEE, 2015, pp. 419–423.
- [16] R. Gallager, "Circularly-symmetric Gaussian random vectors," preprint, 2008.