



**HAL**  
open science

# A Survey on Sound Source Localization in Robotics: from Binaural to Array Processing Methods

Sylvain Argentieri, Patrick Danès, Philippe Souères

► **To cite this version:**

Sylvain Argentieri, Patrick Danès, Philippe Souères. A Survey on Sound Source Localization in Robotics: from Binaural to Array Processing Methods. *Computer Speech and Language*, 2015, 34 (1), pp. 87-112. hal-01058575

**HAL Id: hal-01058575**

**<https://hal.science/hal-01058575>**

Submitted on 27 Aug 2014

**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.

# A Survey on Sound Source Localization in Robotics: from Binaural to Array Processing Methods

S. Argentieri<sup>a,b</sup>, P. Danès<sup>c,d</sup>, P. Souères<sup>c</sup>

<sup>a</sup>*Sorbonne Universités, UPMC Univ. Paris 06, UMR 7222, ISIR, F-75005 Paris, France*

<sup>b</sup>*CNRS, UMR 7222, ISIR, F-75005 Paris, France*

<sup>c</sup>*CNRS, LAAS, 7 avenue du colonel Roche, F-31400 Toulouse, France*

<sup>d</sup>*Univ. de Toulouse, UPS, LAAS, F-31400 Toulouse, France*

---

## Abstract

This paper attempts to provide a state-of-the-art of sound source localization in Robotics. Noticeably, this context raises original constraints—e.g. embeddability, real time, broadband environments, noise and reverberation—which are seldom simultaneously taken into account in Acoustics or Signal Processing. A comprehensive review is proposed of recent robotics achievements, be they binaural or rooted in Array Processing techniques. The connections are highlighted with the underlying theory as well as with elements of physiology and neurology of human hearing.

*Keywords:* Robot audition, source localization, binaural audition, array processing

---

## 1. Introduction

“Blindness separate us from things but deafness from people” said Helen Keller, a famous American author who was the first deafblind person to obtain a Bachelor in Arts, in 1904. Indeed, hearing is a prominent sense for communication and socialization. In contrast to vision, our perception of sound is nearly omnidirectional and independent of the lighting conditions. Similarly, we are able to process sounds issued from a nearby room without any visual information on their origin. But human capabilities are not limited to sound *localization*. We can also *extract*, within a group of speakers talking simultaneously, the utterance emitted by the person we wish to focus on. Known as the term *Cocktail Party Effect* [1], this separation capacity enables us to process efficiently and selectively the whole acoustic data coming from our daily environment. Sensitive to the slightest tone and level variations of an audio message, we have developed a faculty to *recognize* its origin (ringtone, voice of a colleague, etc.) and to *interpret* its contents. All these properties of localization, extraction, recognition and interpretation allow us to operate in dynamic environments, where it would be difficult to do without auditory

---

\*This work is conducted within the European FP7 TWO!EARS project under grant agreement n°618075.

*Email addresses:* [sylvain.argentieri@upmc.fr](mailto:sylvain.argentieri@upmc.fr) (S. Argentieri), [patrick.danes@laas.fr](mailto:patrick.danes@laas.fr) (P. Danès), [philippe.soueres@laas.fr](mailto:philippe.soueres@laas.fr) (P. Souères)

information. All the above impressive Human capabilities have stimulated developments in the area of *Robot Audition*. Likewise, the recent research topic of Human-Robot Interaction may have constituted an additional motivation to investigate this new field, with the aim to artificially reproduce the aforementioned localization, extraction, recognition and interpretation capabilities. Nevertheless, contrarily to Computer Vision, robot audition has been identified as a scientific topic of its own only since about 15 years. Since then, numerous works have been proposed by a growing community, with contributions ranging from sound source localization and separations in realistic reverberant conditions to speech or speaker recognition in the presence of noise. But as outlined in [2], the robotics context raises several unexpected constraints, seldomly taken into account in Signal Processing or Acoustics. Among them, one can cite:

*Geometry constraint:* Though the aim is to design an artificial auditory system endowed with performances inspired by human hearing, there is no need to restrict the study to a biomimetic sensor endowed with just two microphones. Indeed, bringing redundant information delivered by multiple transducers can improve the analysis and its robustness to noise. Straight connections thus appear with the field of Array Processing. Yet, the robotics context imposes an *embeddability* constraint. While Array Processing can consider large arrays of microphones—e.g. several meters long—, robotics implies a tradeoff between the size of the whole sensor and its performances, so that it can be mounted on a mobile platform. Furthermore, applications in humanoid robotics strongly promote the use of only two microphones.

*Real Time constraint:* Many existing methods to sound analysis rely on heavy computations. For instance, a processing time extending over several tens of seconds is admitted to perform the acoustic analysis of a passenger compartment. Contrarily, localization primitives involved in low-level reflex robotics functions—e.g. sensor-based control or auditive/visioauditive tracking—must be made available within a guaranteed short time interval. So, the algorithms computational complexity is a fundamental concern. This may imply the design of dedicated devices or computational architectures.

*Frequency constraint:* Most sound signals valuable to Robotics are *broadband*, i.e. spread over a wide bandwidth w.r.t. their central frequency. This is the case of voice signals, which show significant energy on the bandwidth [300 Hz–3300 Hz] used for telephony. Consequently, narrowband approaches developed elsewhere do not straightly apply in such broadband contexts. Noticeably, this may imply a higher computational demand.

*Environmental constraint:* Robotics environments are fundamentally dynamic and unpredictable. Contrarily to acoustically fully controlled areas, unexpected noise and reverberations are likely to occur, which depend on the room characteristics —dimensions, walls, type of the building materials, etc.—and may singularly deteriorate the analysis performance. The robot itself participates to these perturbations, because of its self-induced noise, e.g. from fans, motors, and other moving parts. A challenge is to endow embedded sound analysis systems with robustness and/or adaptivity capabilities, able to cope with barge-in situations where both the robot and a human are possibly both speaking together.

Generally, most of embedded auditory systems in robotics follow the following classical bottom-up framework: as a first step, the sensed signals are analyzed to estimate sound sources positions; next the locations are used to separate sound of interests from the sensed mixture in order to provide clean noise or speech signals ; finally, speaker

and speech recognition systems are fed with these extracted signals. Of course, other alternatives have also been proposed [3], but this approach remains by far the most used framework in robot audition. Nevertheless, it exhibits the importance of sound localization in the overall analysis process. It has been indeed the most widely covered topic in the community, and a lot of efforts have been made to provide efficient sound localization algorithms suited to the robotics context. Since, in our opinion, Robot Audition has reached an undeniable level of scientific maturity, we feel that the time has come to summarize and organize the main publications of the literature. This paper then attempts to review the most notable contributions specific to sound source localization. Another intent is to underline their connections with theoretical foundations of the field, including with basics of human physiology and neuroscience.

The paper is organized into two parts. First, binaural methods to sound source localization are reviewed in Section 2, under the angle of performances and human operation. Next, Array processing approaches are expounded in Section 3, with a focus on the specificities raised by the robotics constraints. Finally, a conclusion ends the paper.

## 2. Binaural Approaches to Sound Source Localization in Robotics

This section describes a first set of methods which try to mimic diverse aspects of the human auditory system, thus defining the topic of *binaural robot audition*. The common point to all the following works is the use of only two acoustical sensors, generally positioned inside a human-like pinna. There is an obvious practical interest to develop biomimetic auditory sensors containing a small number of microphones: the size is minimal and the embedded electronics is simplified. Significant advances in understanding the biological processes which enable the handling of acoustic data by humans have been obtained up to the 80s [4]. They constitute the natural basis of binaural developments in robotics. Having this in mind, the successive steps of the sound localization process can be described by following the sound propagation, from the source to the binaural sensor:

- As a first step, a sound wave generated by an external source is modified by the presence of the robotic torso, head and pinnae prior to interact with the ears. The induced scattering and spectral changes must be modeled so as to precisely capture the time-space relationship linking the sound source to the binaural signals. From an engineering point of view, this relationship is captured by the so-called *Head Related Transfer Function* (HRTF), which will be studied in the first subsection.
- Next, human localization capabilities mainly rely on some acoustic features extracted by our ears and integrated by our brain. Those features have been extensively studied; among them, one can cite *Interaural Cues* for horizontal localization, or spectral notches for vertical localization [4]. A lot of them have also been used in robotics. These will be reviewed in the second subsection.
- Finally, on the basis of these features, sound localization itself is performed. Numerous approaches have been proposed so far, and the most prominent one in robotics are outlined in the third subsection.

In all the following, the left and right microphone signals will be referred to as  $l(t)$  and  $r(t)$  respectively, with  $t$  the time (in s). Their frequency counterparts, obtained through a

Fourier Transform, will be denoted by  $L(f)$  and  $R(f)$  respectively, with  $f$  the frequency (in Hz). Sound source position is expressed in terms of horizontal azimuth angle  $\theta$ , elevation angle  $\varphi$  in the median plane, and distance  $r$ , all of them being expressed w.r.t. an origin located at the robot’s head center. In the remaining of the paper, the position  $(\theta, \varphi) = (0^\circ, 0^\circ)$  corresponds to a sound source in front of the head (i.e. at boresight).

## 2.1. The Head Related Transfer Function

### 2.1.1. Definition

The HRTF captures the relationships between the signal  $s(t)$  originating from a sound source and captured at a certain arbitrary reference position in space and the two signals perceived by the two ears. These relationships can be written in the form

$$\begin{cases} L(f) = H_L(r_s, \theta_s, \varphi_s, f)S(f), \\ R(f) = H_R(r_s, \theta_s, \varphi_s, f)S(f), \end{cases} \quad (1)$$

where  $H_L(\cdot)$  and  $H_R(\cdot)$  represent the left and right HRTFs respectively,  $(r_s, \theta_s, \varphi_s)$  is the actual sound source position w.r.t. the chosen reference position, and  $S(f)$  the Fourier transform of  $s(t)$ . Importantly, the HRTFs account for all the modifications brought by the body of the listener, including torso, head and pinnae effects, to the incident sound wave. So it varies significantly from a human listener to another as it reflects his/her own morphology. The same applies in robotics, where all the possible acoustic scatters impact on the sensed signals, and are thus captured by the corresponding HRTFs. But it is fundamental to understand that the HRTF strictly corresponds to a propagation in free field and does not include room reflections nor reverberations. Consequently, the HRTF can be obtained in two ways. The first solution is to accurately model the body and head effects. If the robot shapes are simple, then basic acoustic equations can be sufficient to account for the acoustic effect on the signals. In the case of a more realistic robot, with complex body, shoulders, nose, pinnae, etc., an acoustic simulation software might be necessary to solve the problem through finite-element methods [5]. The second solution consists in identifying the HRTF through real measurements, which must be performed in an anechoic room. This solution may not be so practical for every robotic platform, for it requires specific hardware/software. Hopefully, various databases are proposed in the literature. One can cite, among others, the celebrated CIPIC database [6], published by the CIPIC Interface Laboratory from the University of California Davis, and accessible from the URL <http://interface.cipic.ucdavis.edu/>, or the recent HRTF database proposed by the Deutsche Telekom Laboratories (TU-Berlin) [7], located at <https://dev.qu.tu-berlin.de/projects/measurements/wiki>

Typical HRTFs extracted from the CIPIC database is shown on Fig. 1, for the azimuth  $\theta = 35^\circ$  and elevation  $\phi = 20^\circ$ . Its time counterpart Head Related Impulse Response (HRIR) is also represented. This figure highlights the famous shadowing effect of the head: depending on the source position, the left and right signals differ in terms of time of arrival (cf. the delay between the left and right HRIR), but also in terms of spectral content (cf. the amplitude difference between the two HRTFs and the spectral notches positions). These last cues will often serve as the basis to infer localization (see §2.2). Readers interested in a more complete tutorial on HRTF can refer to [8], where experimental and theoretical data are compared.

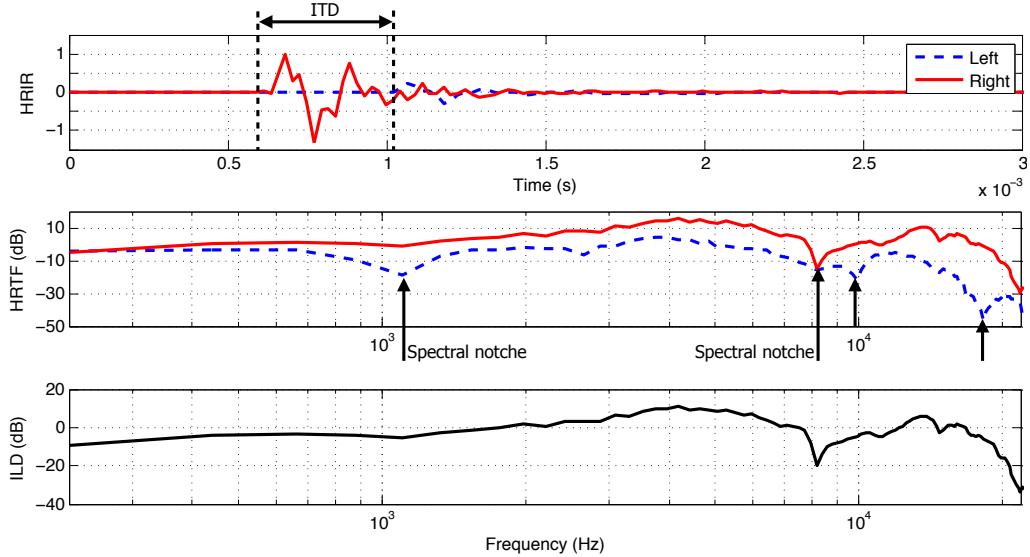


Figure 1: HRIR and HRTF data for a subject of the CIPIC database [6]. Interaural (ITD/ILD) and monaural (spectral notches) cues are also reported.

### 2.1.2. HRTF models in Robotics

As already outlined, the complex structure of most robotic platforms prevents the access (through simulations or identifications) to the exact robot HRTFs at the two ears. Consequently, some simple head models have been proposed by the robot audition community, with the aim to capture the robot head effect on the binaural signals up to some extent.. The three most classical models are depicted on Fig. 2. They consist in considering the left and right microphones in the free field —Auditory Epipolar Geometry (AEG) [9]—, placed at the surface of a disk —Revised Epipolar Geometry (RAEG) [10]—, or at the surface of a sphere —Scattering Theory (ST) [11]—. AEG and RAEG are the most elementary model. Provided that  $\theta$  and  $f$  stand for the azimuth and frequency of a farfield source, the (R)AEG approximations of the left and right HRTFs  $H_L^{(R)AEG}(\cdot)$  and  $H_R^{(R)AEG}(\cdot)$  write as

$$\begin{cases} H_L^{(R)AEG}(\theta, f) = 1, \\ H_R^{(R)AEG}(\theta, f) = e^{-j\phi(\theta)} = e^{-j2\pi f\tau_{(R)AEG}(\theta)}, \end{cases} \quad (2)$$

highlighting the fact that the two binaural signals only differ by a delay  $\tau_{(R)AEG}(\theta)$  which is a function of the source angle (note that the left channel has been arbitrarily considered here as the reference). The third ST (spherical) model is more involved. Let  $\beta$  be the so-called incidence angle, i.e. the angle between the line from the center of the sphere approximating the head to the source position  $(r_s, \theta_s)$ , and the line from the center of the same head to a measurement point at which the HRTF must be computed. Considering a perfect rigid spherical head, the expression of the diffracted sound pressure

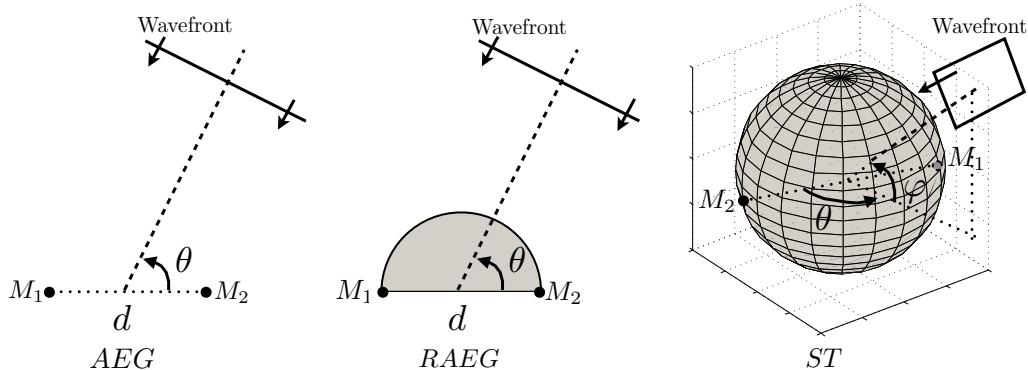


Figure 2: The three classical head models: auditory epipolar geometry (AEG, left), revised auditory epipolar geometry (RAEG, middle), and scattering theory (ST, right).

wave received at the measurement point allows to write [12]:

$$H^{\text{ST}}(r, \beta, f) = \frac{rce^{-jr2\pi f/c}}{ja^22\pi f} \sum_{m=0}^{\infty} (2m+1)P_m[\cos(\beta)] \frac{h_m(r2\pi f/c)}{h'_m(a2\pi f/c)}, \quad (3)$$

where  $H^{\text{ST}}(r, \beta, f)$  is the transfer function linking the sound pressure received at the measurement point and the free-field pressure existing at the head center, with  $c$  the speed of sound and  $a$  the head radius.  $P_m(\cdot)$  and  $h_m(\cdot)$  are the Legendre polynomial of degree  $m$  and the  $m^{\text{th}}$ -order spherical Hankel functions respectively, while  $h'_m(\cdot)$  denotes the derivative of the function  $h_m(\cdot)$ . Assuming that the two microphones are antipodally placed on the surface of the sphere, the left and right HRTFs, respectively denoted by  $H_L^{\text{ST}}(r, \theta, f)$  and  $H_R^{\text{ST}}(r, \theta, f)$ , are then given by, for a sound source located at  $(r, \theta)$ ,

$$\begin{cases} H_L^{\text{ST}}(r, \theta, f) = H^{\text{ST}}\left(r, -\frac{\pi}{2} - \theta, f\right), \\ H_R^{\text{ST}}(r, \theta, f) = H^{\text{ST}}\left(r, \frac{\pi}{2} - \theta, f\right). \end{cases} \quad (4)$$

## 2.2. Binaural and monaural cues for localization

Once the link between the sound source signal to be localized and the two resulting binaural signals has been modeled, it is necessary to focus on the binaural features which can be extracted from these signals to infer localization. These features are first recalled through a short review on sound source localization in humans. Then, the way how these cues can be coupled with the aforementioned HRTFs is investigated.

### 2.2.1. Sound source localization in humans

About 100 years ago, Rayleigh proposed the *duplex theory* [13], which explains that horizontal localization is mainly performed through two primary auditory cues, namely the *Interaural Level Difference* (ILD) and the *Interaural Time Difference* (ITD). The ILD relates to the difference between the intensity of signals perceived by the right and left ears, due to the head frequency-dependent scattering. Noticeably, if a source emits at

a frequency higher than about 750 Hz, then the head and any small-sized element of the face induce scattering, which significantly modifies the perceived acoustic levels, so that the ILD can exceed 30 dB. On the contrary, the ILD is close to 0 dB at low frequencies, as fields whose wavelengths are greater than the head diameter undergo no scattering. This property can clearly be deduced from Fig. 1, where the left and right HRTF amplitudes only significantly differ for frequencies greater than about 800 Hz. The second auditory cue is known as the *Interaural Phase Difference (IPD)*—or its time-counterpart termed *Interaural Time Difference (ITD)*. The ITD is justified by the path difference to be traveled by the wave to reach the ears. It appears on Fig. 1 as a delay between the two HRIR onsets. Note however that the maximum value involved in localization is around  $700\mu\text{s}$ —i.e. one period of a 1400 Hz sound— due to the ambiguity of IPD values greater than  $2\pi$ . So, two frequency domains can be exhibited in human horizontal localization, each one involving a distinct acoustic cue. Frequencies under  $\sim 1$  kHz are azimuthally localized by means of the IPD, while frequencies above  $\sim 3$  kHz exploit the ILD.

It can be straightly inferred that a source emitting from the vertical symmetry plane of the head produces no interaural difference. However, Humans are still able to perform a localization in such conditions. Indeed, obstacles—including shoulders, head, outer ear, etc.—play the role of scatterers which modify the frequency components of acoustic waves. This filtering effect is essential in our ability to determine the elevation of a sound source. Indeed, the sum of all the reflections occurring around the head induces notches into the perceived spectrum, the positions of which are significantly affected by the source elevation [14], see Fig. 1. The acoustic feature for vertical localization is thus a spectral cue, termed "*monaural*" as it involves no comparison between the signals perceived at the two ears. Consequently, these notch positions are likely to be used by the brain to infer the elevation.

While the directional aspects of localization have been widely studied, distance perception has received substantially less attention. Generally, it is admitted that like angular estimations, sound source distance can also be inferred from various acoustical properties of the sounds reaching the two ears. Known distance dependent acoustical cues include sound intensity, which unfortunately also depends on the intrinsic source energy, as well as on interaural differences, on the spectral shape and on the Direct-to-Reverberant sound energy Ratio (DRR). So, one can see that nearly all the aforementioned cues, which are used to estimate the angular position of a sound source, are also directly linked to the distance parameter. Actually, human most likely combine these distant-dependent cues together with *a priori* information on the surrounding space so as to get the sensation of a stable distance. The reader interested in this topic will find a comprehensive review in [15]. But one has to keep in mind that Human performances in distance discrimination are quite poor: listening tests have proven that humans use to significantly overestimate the distance to sources closer than 1m, while we underestimate distances greater than 1m [15].

### 2.2.2. Auditory models in Robotics

As shown in the previous paragraph, the head effect on the perceived sounds is frequency dependent. Such a frequency decomposition of the signals is often implemented with a FFT, while other authors proposed the use of classical bandpass filters, see [16]. Nevertheless, how efficient they may be, these methods do not perform a frequency decomposition similar to the human inner ear. *Gammatone filters*, modeling the vibra-



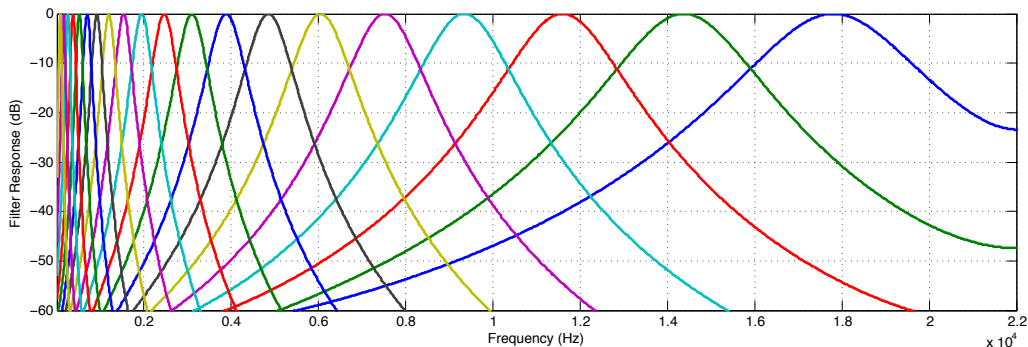


Figure 3: Typical Gammatone filters frequency responses.

tion of the basilar membrane inside the cochlea, have proven well suited to describe the cochlear impulse response in the cat's cochlea [17]. They also constitute a good approximation of human spectral analysis at moderate sound levels [18]. Typical gammatone filters frequency responses are reported in Fig. 3. It can be seen that their bandwidths increase with frequency, in such a way that they represent around 15% of the center frequencies. This is one of the main features of the human auditory system, which confirms our better ability to discriminate low frequencies [19]. As a consequence of this decomposition, the representation of any sound signal information is more likely close to the human perception of sounds. As it will be shown in the following, this gammatone frequency decomposition is now very commonly used. Readers interested in more involved auditory models can refer to the Auditory Modeling Toolbox [20], available at <http://amtoolbox.sourceforge.net>

### 2.2.3. Binaural cues for horizontal localization in Robotics

Among all the acoustic features that can be extracted with two microphones, binaural cues are the most often used in robotics. The IPD and/or ILD can indeed lead, with just two microphones and simple computations, to a localization in azimuth. Let  $IPD_{\text{exp}}(f)$  and  $ILD_{\text{exp}}(f)$  term the experimental IPD and ILD extracted from the two signals. There exist numerous way to estimate these IPD and ILD values: computations in the time or frequency domain, bioinspired models, etc. Readers interested in a review of these approaches can consult [16]. Whatever the approach, from these experimental values, the problem is then to determine the position of the emitting sources. This involves a model, expressed either in a mathematical closed-form or as experimental relationships uniting the source attributes (position, frequencies,...) and the induced auditory cues.

Considering the AEG or RAEG model represented by Eq. (2), the delay  $\tau_{(R)AEG}(\theta)$  represents the ITD, and its phase counterpart, i.e. the IPD, then verifies [10]

$$\begin{cases} IPD_{AEG}(\theta, f) = 2\pi f \tau_{AEG}(\theta) = \frac{2\pi f a}{c} \cos \theta, & (5a) \\ IPD_{RAEG}(\theta, f) = 2\pi f \tau_{RAEG}(\theta) = \frac{\pi f a}{c} \left( \left( \frac{\pi}{2} - \theta \right) + \cos \theta \right). & (5b) \end{cases}$$

The main advantage of the first AEG formulation is that the azimuth  $\theta$  can straightly be

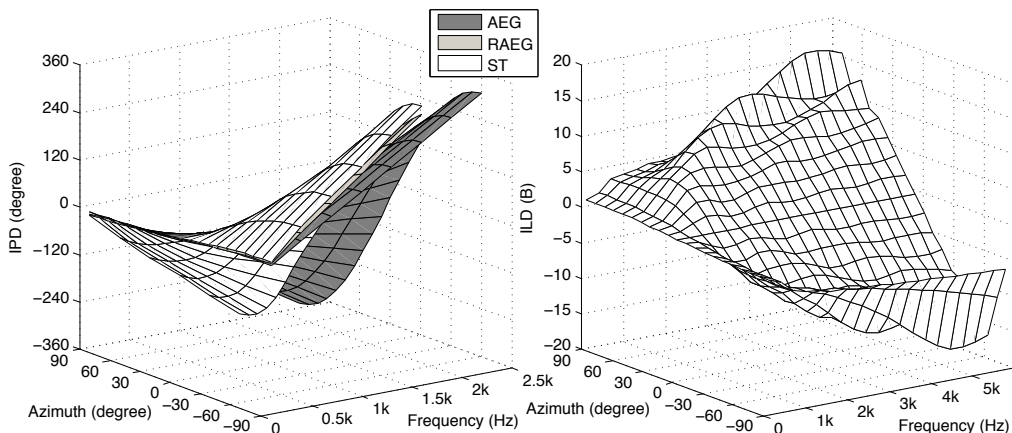


Figure 4: Comparison of the AEG, RAEG and ST models. (Left) IPD as a function of azimuth  $\theta$  and frequency. (Right) ILD values for the ST model.

approximated by inverting (5a), assimilating  $IPD_{AEG}(\theta, f)$  to the experimental  $IPD_{exp}$ . However, as already outlined, it cannot describe the effect of a head located between the two microphones, inducing scattering of the sound wave. To better take into account its presence, the RAEG model can be used (note that RAEG is analog to the classical Woodworth-Schlosberg formalization [21]). Indeed, results from [10] show that simulations obtained from this model fit experimental measurements in an anechoic room within the frequency band [500 Hz–800 Hz]. But the RAEG model, while being more realistic than AEG, does not fully account for the head waveguide effect. Additionally, both do not provide any meaningful value for the ILD cue, which is accounted for in the theoretical spherical model. Indeed, one has for this ST model

$$\begin{cases} IPD_{ST}(r, \theta, f) = \arg(H_L^{ST}(r, \theta, f)) - \arg(H_R^{ST}(r, \theta, f)), & (6a) \\ ILD_{ST}(r, \theta, f) = 20 \log_{10} \frac{|H_L^{ST}(r, \theta, f)|}{|H_R^{ST}(r, \theta, f)|}. & (6b) \end{cases}$$

Compared to epipolar geometries, the scattering theory exhibits more reliable theoretical expressions of both IPD and ILD as functions of the azimuth  $\theta$  and distance  $r$ , and can thus lead to more reliable identification schemes for localization. It is however important to notice that the accuracy of the approach still depends on the capacity to cope with the room acoustics, which is not always possible in practice. Indeed, if the binaural cues are obtained inside in a real robotics environment including noisy sound sources and reverberation, the results may get very bad: since the models do not capture the distortion due to the room acoustics, the theoretical binaural cues and the measured one cannot fit with each. Nevertheless, the ST formalism was exploited in [11] and [22] to express the pressure field perceived by two microphones symmetrically laid over a sphere, and experimentally tested by Handzel *et al.* [23] on a spherical plastic head. But whatever the model, and as outlined in Human audition, the observed inappropriateness of IPD (resp. ILD) for high (resp. low) frequencies extends outside the scope of AEG, RAEG

and ST strategies, as can be seen in Fig. 4. As mentioned in §2.2.1, frequencies above about 1400Hz lead to IPD values greater than  $2\pi$  and becomes ambiguous. Noticeably, the human auditory system then relies on the ILD at these frequencies. Indeed Fig. 4 exhibits high ILD values, reaching up to 25dB for this frequency domain, in the ST model.

### 2.3. Exploitation in Robotics

Historically, most initial contributions to robot audition were rooted into the binaural paradigm. However, as shown in the following, the results remained mixed when facing real-life environments, involving noises and reverberations together with wideband non-stationary sources.

#### 2.3.1. Horizontal localization

In the early 2000s, the use of interaural difference functions for azimuth localization was deeply studied in the framework of the SIG project [9, 24]. As the robot cover cannot be perfectly isolated from internal sounds, an adaptive filter exploiting the data provided by inner microphones was used to eliminate the motor noises (e.g. ego-noise mentioned in section 1) from the audio signal perceived by the external pair of microphones. This active auditory system thus allowed to perform measurement during motion [9], and constituted an interesting improvement over former methods—for instance [25] on the COG humanoid, or [26]—based on the *stop-perceive-act* principle. On this basis, an “Active Direction Pass Filter” (ADPF) grounded on the ST model was proposed in [27] to determine the origin of a sound source and extract it out of a mixture of surrounding sounds. This *model-matching* approach has since been used in a lot of contributions. For a fixed distance  $r_s$ , for all frequencies under (resp. above)  $f_{th} = 1500Hz$ , and for each  $\theta$ , the system computes the theoretical  $IPD_{ST}(r_s, \theta, f)$  (resp.  $ILD_{ST}(r_s, \theta, f)$ ) through (6). Then a cost function  $d_{IPD}$  (resp.  $d_{ILD}$ ) is defined to measure the distance between the measured  $ITD_{exp}(f)$  (resp.  $ILD_{exp}(f)$ ) interaural differences and theoretical ones. The two distances  $d_{IPD}$  and  $d_{ILD}$  are then integrated into a belief factor  $P_{IPD+ILD}(\theta)$ . The angle  $\hat{\theta}_s = \arg \max_{\theta} P_{IPD+ILD}(\theta)$  is then regarded as the sound source azimuth. The sound source separation performances were also evaluated and compared in [28], depending on the model (RAEG *vs* ST) that relates the source azimuth and the interaural cues. Clearly, the scattering theory provided the best results. So far, these contributions have been among the rare complete binaural audition systems, integrating localization, source separation and recognition. Nevertheless, since HRTFs do not capture the acoustic response of the room where the robot operates, their applicability is generally limited to well-identified environments. One solution could consist in learning the head effect in realistic conditions. Such an idea was successfully assessed in [29] through a dedicated neural network able to generalize learning to new acoustic conditions. One can also cite [30], or [31], where the iCub humanoid robot’s head was endowed with two pinnae. The localization is performed by mapping the aforementioned sound features to the corresponding location of the source through a learning method. Another approach is proposed in [32]. Auditory events corresponding to relevant ITD values are gathered into histograms, which are then approximated by Gaussian models whose parameters are identified through the EM method. Peaks in the resulting histogram are then regarded as potential sound azimuths. This allows to cope with the multisource case, where

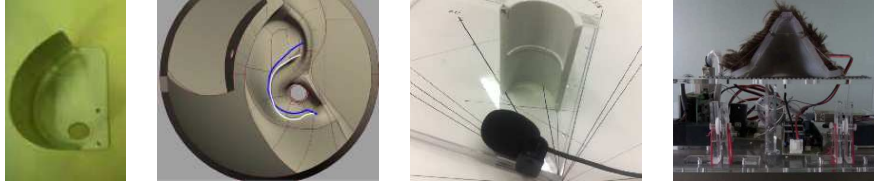


Figure 5: Some artificial pinnas from the literature. Pictures extracted from (left to right) : [37], [35], [39], [40]. They all share the fundamental asymmetry property.

multiple sound sources are likely active at the same time. Finally, an implementation of the celebrated biology-inspired Jeffress model is proposed in [33] on a simple robot head endowed with two microphones and stereovision. Interestingly, the ILD pathway is also modeled with a 2D spiking map. The merging of the two interaural maps is also addressed, so as to obtain an efficient sound localization system. The proposed method is shown to share some common well-known properties of the human auditory system, like the ITD maximal efficiency reached when the sound source is in front of the observer. But whatever the approach, ITDs and ILDs can be extracted from the binaural signals in numerous ways: through correlation [34], zero-crossing times comparison [35], or in the spectral domain [36]. A systematic study of binaural cues, and an analysis of their robustness w.r.t. reverberations is proposed in [16]. Results show that binaural cues extracted from gammatone filters outperforms other techniques.

### 2.3.2. Vertical localization: spectral cues

As indicated in §2.2.1, the elevation of a sound source is mainly related to the positions of notches in the spectra of the perceived signals, which stem from acoustic reflections due to the head and the outer ear. In robotics, quite few authors have developed techniques based on spectral cues. Most of them are based on the scattering induced by an artificial pinnae in charge of collecting the acoustic pressure information and of driving it to microphones. For humans, the specific shape of the pinnae enables a selective spatial amplification of acoustic pressure variations, with a quality factor reaching up to 20 dB. Reproducing such capabilities in Robotics is a difficult problem due to the lack of a model of the pinnae shapes which lead to elevation dependent notches. Yet, as a rule of thumb, these shapes must be irregular or asymmetric, and artificial pinnae were proposed in [37], [38], [31], [35], or [39]. Fig. 5 shows some of them. A simplified model, inspired by [41] and based on the superposition of the incident wave with a single wave reflected by the pinnae, enables the prediction of the elevation from the position of notches. Noticeably, these notches, which appear or disappear depending on the elevation, may be hard to detect or may even be blurred by spurious notches induced by destructive interferences coming from acoustic reflections on obstacles. To solve this problem, [31] introduces the *interaural spectral difference* as the ratio between the left and right channel spectra. While notches may be indistinct in the complex spectra of the two signals, the interaural spectral difference, when interpolated with a 12-degree polynomial, can enable the extraction of their frequency positions. Another solution is proposed in [35]. It consists in computing the difference of the left and right energies coming from 100 frequency channels, ranging from 100Hz to 20kHz. Strictly speaking, this approach does

not involve monaural cues anymore, but allows to obtain spectral cues which are said less sensitive to the source signal frequency content. Concerning the design of the pinnae, a model including a more extensive description of the reflected and diffracted sound waves is proposed in [42]. Though it leads to new theoretical expressions of the spectra, it remains hardly valuable for the design of artificial outer ears. Reference [39] also exhibits four different pinnae together with their induced frequency responses for an original work on sound localization from a single microphone. On the other hand, inspired by animals that are able to change the configuration of their pinnae, [40] proposed an *active* ear, which is able to modify its shape to encode elevation (and azimuth).

### 2.3.3. Distance localization

In the topic of robot audition, distance estimation has been so far based on the triangulation idea. For instance, on the basis of the estimation  $\hat{\theta}_1$  and  $\hat{\theta}_2$  of the azimuth at two distinct positions, triangulation allows to estimate the distance between the robot and the sound source, together with the source azimuth. Generally, this is only possible if the sound source is static in the environment. Some recent works [43, 44] proposed a filtering strategy to cope with a possibly moving source. The algorithm mainly relies on ITD to provide an estimation of the source position  $(r, \theta)$  during the movement of a binaural sensor. Distance estimation is also investigated in [45]. Several auditory cues, like interaural differences, sound amplitude and spectral characteristics are compared. Convincing results are shown, exhibiting an estimation error lower than 1m for a 6m-far sound source. But the author outlines that its study doesnot capture the full variability of natural environments. Recent contributions also propose to estimate the DRR. Indeed, it has been shown that distance estimation by humans are more accurate in a reverberant space than in an anechoic one. This estimation is not straightforward: [46] proposes a binaural equalization-cancellation technique, while [47] hypothesises the use of the frequency dependent magnitude squared coherence between the left and right signals.

### 2.4. Conclusion

Using very few microphones, interesting developments have been proposed by the Robotics community to provide the robots with a first ability to localize sound sources in their environment. The results are however contrasted. Reproducing the auditory faculty of the human ear is a very difficult problem. First, the exploitation of interaural cues requires a very precise modeling of the perturbations induced by the presence of the head. Second, binaural cues are still very hard to exploit. In any case, an accurate model of the propagation turns out to be essential to finely describe the evolution of auditory cues. Furthermore, all these techniques appear to be very sensitive to variations of the acoustic environment. Most models have been experimentally validated in an anechoic room but cannot be used to accurately localize sounds in real conditions, unless a precise description of the robot's environment is given. But recent *active* variations of the existing algorithms have recently benefited from the additional information brought by the robot motion, thus renewing the interest in binaural approaches to sound localization. Nevertheless, all these difficulties have motivated the Robotics community to also envisage localization methods based on an higher number of microphones, possibly benefiting from existing signal processing advances. An overview of these techniques is presented in the next section.

### 3. Array processing approaches to localization in Robot Audition

This section deals with the second paradigm mainly used in robot audition: microphone arrays. Contrarily to binaural approaches, where only two sensors are used, array processing relies on multiple microphones, spatially organized along various geometries (such as a line, a circle, a sphere, or the vertices of a cube). Thanks to the redundancy in the signals from the various channels, the acoustic analysis performance and/or robustness can be improved [48]. Multiple contributions have been proposed in a robotics context, generally concerning source detection and localization, source separation, and speaker/speech recognition. Again, this section is entirely devoted to sound source localization, by focusing only on methods used in robotics.

After having introduced some notations, the celebrated signal processing method MUSIC (Multiple Signal Classification) is hereafter presented. This method, though very powerful, exemplifies the limits imposed by the real time constraint. Next, approaches relying on the temporal delays due to the wave propagation between the microphones are overviewed. They illustrate how information redundancy can enhance the localization accuracy and robustness. The section ends with beamforming-based methods. Their simplicity and ease of implementation makes them ideal candidates for an application in robotics.

#### 3.1. Theoretical aspects of array processing in Robotics

##### 3.1.1. Notations and definitions

Consider  $S$  pointwise sound sources emitting at locations referenced by  $\mathbf{r}_s^s = (r_s, \theta_s, \varphi_s)$ ,  $s = 1, \dots, S$ , in a spherical coordinates system. In the following, any monochromatic space-time signal reads as  $y(\mathbf{r}, t) = Y(\mathbf{r}, k)e^{jkct}$ , with  $k = 2\pi f/c$  the wavenumber. In addition, let a microphone array be composed of  $N$  identical omnidirectional microphones placed at locations  $\mathbf{r}_n^m$ ,  $n = 1, \dots, N$ . Then, the sound signal  $m_n(t)$  issued by the  $S$  sources and perceived by the  $n^{\text{th}}$  transducer can be written as

$$m_n(t) = \sum_{s=1}^S \frac{\|\mathbf{r}_s^s\|}{\|\mathbf{r}_n^m - \mathbf{r}_s^s\|} s_s^0\left(t - \frac{\|\mathbf{r}_n^m - \mathbf{r}_s^s\|}{c} + \frac{\|\mathbf{r}_s^s\|}{c}\right) + b_n(t), \quad (7)$$

where  $s_s^0(t)$  terms the fictitious signal perceived at  $\mathbf{r} = \mathbf{0}$  and stemming from the single  $s^{\text{th}}$  source, and the additive noise  $b_n(t)$  accounts for parasitic sources in the environment as well as electronic noise in the microphones outputs. So, the Fourier transforms  $M_n(k)$ ,  $S_s^0(k)$ ,  $B_n(k)$  of  $m_n(t)$ ,  $s_s^0(t)$ ,  $b_n(t)$  satisfy

$$M_n(k) = \sum_{s=1}^S V_n(\mathbf{r}_s^s, k) S_s^0(k) + B_n(k), \quad (8)$$

with

$$V_n(\mathbf{r}, k) = \|\mathbf{r}\| e^{jk\|\mathbf{r}\|} \frac{e^{-jk\|\mathbf{r}_n^m - \mathbf{r}\|}}{\|\mathbf{r}_n^m - \mathbf{r}\|} \quad (9)$$

the  $n^{\text{th}}$  entry of the *array*—or *steering*—*vector*  $\mathbf{V}(\mathbf{r}, k) \triangleq (V_1(\mathbf{r}, k), \dots, V_N(\mathbf{r}, k))^T$ . Defining the *source*, *observation* and *noise vectors* by  $\mathbf{S}^0(k) \triangleq (S_1^0(k), \dots, S_S^0(k))^T$ ,

$\mathbf{M}(k) \triangleq (M_1(k), \dots, M_N(k))^T$  and  $\mathbf{B}(k) \triangleq (B_1(k), \dots, B_N(k))^T$  respectively, (8) can be turned into the matrix form

$$\mathbf{M}(k) = \mathcal{V}(\mathbf{r}_1^s, \dots, \mathbf{r}_S^s, k) \mathbf{S}^0(k) + \mathbf{B}(k), \quad (10)$$

with  $\mathcal{V}(\mathbf{r}_1^s, \dots, \mathbf{r}_S^s, k) \triangleq (\mathbf{V}(\mathbf{r}_1^s, k) | \dots | \mathbf{V}(\mathbf{r}_S^s, k))$  the *array matrix*. Note that (9) can significantly be simplified when the distance to the sources tends to infinity, as the wavefronts become planar. This simplification defines the “farfield hypothesis”. In the following, quantities related to farfield will be superscripted by the symbol  $\infty$ , so that the *farfield array vector* writes as  $\mathbf{V}^\infty(\theta, \varphi, k) \triangleq (V_1^\infty(\theta, \varphi, k), \dots, V_N^\infty(\theta, \varphi, k))^T$ , with  $V_n^\infty(\theta, \varphi, k) = V_n^\infty(\mathbf{r}, k) \triangleq \lim_{r \rightarrow \infty} V_n(\mathbf{r}, k)$ .

From now on, let's consider a linear microphone array, constituted of  $N$  microphones located at  $z_1, \dots, z_N$  along the  $\mathcal{Z}$ -axis. Consequently, because of the rotational symmetry of the problem, all characteristics are invariant w.r.t. the elevation  $\varphi$ , so that the location vector  $\mathbf{r} = (r, \theta, \varphi)$  reduces to  $\mathbf{r} = (r, \theta)$ . In addition, the  $n^{\text{th}}$  entry (9) of the array vector becomes

$$V_n(\mathbf{r}, k) = V_n(r, \theta, k) = \frac{r e^{jkr} e^{-jk\sqrt{r^2 + z_n^2 - 2rz_n \cos \theta}}}{\sqrt{r^2 + z_n^2 - 2rz_n \cos \theta}}. \quad (11)$$

In the farfield, (11) particularizes into the well-known expression  $V_n^\infty(\theta, k) = e^{-jkz_n \cos \theta}$ .

### 3.1.2. The MUSIC method

The MUSIC (MULTiple Signal Classification) method, initially proposed in [49], belongs to the so-called “high resolution” approaches because of the sharpness of the conclusions it provides. It is so far one of the most used algorithm in Robotics. The pointwise sound sources to be localized are assumed independent, zero-mean stationary, of single frequency  $k_0$ , and in number  $S < N$ . In equation (10), the additive noise is assumed zero-mean, stationary, temporally and spatially white, of known equal power on each microphone, and independent of the sources. So, denoting by  $\mathcal{I}$  and  $\mathcal{O}$  the identity and zero matrices, and  $E[\cdot]$  the expectation operator, it is supposed that

$$\mathbf{\Gamma}_B = E[\mathbf{B}\mathbf{B}^H] = \sigma_N^2 \mathcal{I}_{N \times N} \text{ and } E[\mathbf{S}^0 \mathbf{B}^H] = \mathcal{O}_{S \times N}. \quad (12)$$

As has just been done, the dependencies of variables upon the single involved wavenumber  $k_0$  will be temporarily omitted. MUSIC determines the sources number  $S$  together with their ranges and azimuths from the eigendecomposition of the covariance—or interspectral— $N \times N$  matrix  $\mathbf{\Gamma}_M = E[\mathbf{M}\mathbf{M}^H]$  relative to the signals perceived at the array.

$$\mathbf{\Gamma}_M = (\mathbf{u}_S | \mathbf{u}_N) \begin{pmatrix} \lambda_1 + \sigma_N^2 & & \mathcal{O} & | & \\ & \ddots & & | & \mathcal{O} \\ & & \lambda_S + \sigma_N^2 & | & \\ - & \mathcal{O} & - & | & \sigma_N^2 \bar{\mathcal{I}}_{N-S} \end{pmatrix} (\mathbf{u}_S | \mathbf{u}_N)^H, \quad (13)$$

where the real  $\lambda_1, \dots, \lambda_S$  are sorted increasingly  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_S > 0$ ,  $\mathbf{u}_S = (\mathbf{u}_1 | \dots | \mathbf{u}_S) \in \mathbb{C}^{N \times S}$  and  $\mathbf{u}_N = (\mathbf{u}_{S+1} | \dots | \mathbf{u}_N) \in \mathbb{C}^{N \times (N-S)}$ . The right eigenvectors  $\mathbf{u}_1, \dots, \mathbf{u}_S$  related to the  $S$  greatest eigenvalues  $\lambda_1 + \sigma_N^2, \dots, \lambda_S + \sigma_N^2$  of  $\mathbf{\Gamma}_M$

can be shown to span the range of  $\mathcal{V}(\mathbf{r}_1^s, \dots, \mathbf{r}_S^s)$ , i.e. the  $S$ -dimensional subspace  $\mathcal{S}$  of  $\mathbb{C}^N$  generated by the steering vector evaluated at the sources locations, henceforth termed *signal space*. In the same way, the range of the matrix  $\mathcal{U}_N$  of the  $N - S$  remaining eigenvectors—associated to the eigenvalues  $\sigma_N^2$ —is henceforth termed the *noise space*  $\mathcal{N}$ . Noticeably, the full eigenvectors matrix  $(u_s | u_N)$  can be selected as orthogonal, i.e.  $(u_s | u_N)^H (u_s | u_N) = \mathcal{I}_N$ . Consequently, under the aforementioned statistical hypotheses, (13) enables the recovery, from the covariance matrix  $\mathbf{\Gamma}_M$ , of the number of sources—which is  $N$  minus the number of repetitions of  $\sigma_N^2$ —and of their locations—for their associated steering vectors are orthogonal to  $\mathcal{U}_N$ —. But in practice,  $\mathbf{\Gamma}_M$  is not known, as only one time record of  $\mathbf{m}(t) \triangleq (m_1(t), \dots, m_N(t))^T$  is available. One common strategy consists in computing an approximation of this quantity on  $W$  time snapshots, e.g. by defining

$$\widehat{\mathbf{\Gamma}}_M = \frac{1}{W} \sum_{w=0}^{W-1} \widehat{\mathbf{M}}_w(k) \widehat{\mathbf{M}}_w^H(k), \quad (14)$$

where  $\widehat{\mathbf{M}}_w(k)$  denotes an approximation of  $\mathbf{M}(k)$  from a  $L$ -point Discrete Fourier Transform (DFT) on the  $w^{\text{th}}$  time snapshot. Finally, the locations of the sound sources are established by isolating the maximum values of the *pseudo-spectrum*

$$h(r, \theta) = \frac{1}{\mathbf{V}^H(r, \theta) \widehat{\Pi}_N \mathbf{V}(r, \theta)}, \quad (15)$$

where  $\widehat{\Pi}_N = \widehat{\mathcal{U}}_N \widehat{\mathcal{U}}_N^H$  is called the *projector onto the noise space* and is estimated through the eigendecomposition of  $\widehat{\mathbf{\Gamma}}_M$ . All these developments have been obtained for a single frequency  $k_0$ . Since most sources of interest in robotics are not narrowband, broadband extensions must be proposed to cope with realistic scenarios. These will be mentioned in §3.2.1.

### 3.1.3. Localization through correlation

In the same way as a sound reaching our two ears is delayed due to propagation, the spatial sampling performed by a microphone array induces temporal delays, also termed *Time Delay(s) of Arrival*, or TDOAs. The approaches outlined in this section aim at estimating the delay  $\Delta T_{ij}$  between a pair  $i, j$  of microphones constituting the array through the computation of a correlation function  $R_{m_i m_j}$ . Noticeably, the notions of TDOA and of ITD/IPD—see §2.2.1—are fairly similar. Despite ITDs/IPDs are generally devoted to binaural approaches, both quantities account for the same physical reality, and can then be estimated in the same way. Furthermore, some biological models claim that the ITD/IPD interaural cue is determined by the brain through a correlation involving dedicated neuronal delay lines [50], sometimes called *coincidence detectors* [51]. Nevertheless, as the forthcoming TDOA computations as well as their exploitation significantly differ from the functioning of the brain, they have been classified into the array processing approaches.

In all this subsection, the link between the two signals measured on the  $i^{\text{th}}$  and  $j^{\text{th}}$  microphone of the array (with  $i \neq j$  in all the following) is modeled along

$$\begin{cases} m_i(t) = s(t) + n_i(t) \\ m_j(t) = (s * h_{\mathbf{r}})(t) + n_j(t), \end{cases} \quad (16)$$



with  $*$  the convolution operator,  $s(t)$  the signal received on the  $i^{\text{th}}$  arbitrarily chosen microphone and originating from the source to be localized, and  $h_{\mathbf{r}}(t)$  the deterministic impulse response between the two considered signals.  $s(t)$ ,  $n_i(t)$  and  $n_j(t)$  are also hypothesized as zero-mean stationary signals. If the signal  $s(t)$  propagates from the source to the array in the free field, and without any scatters placed in the vicinity of the microphones, the impulse response  $h_{\mathbf{r}}(t)$  only captures the TDOA  $\Delta T_{ij}$  between the two receivers, i.e.  $h_{\mathbf{r}}(t) = \delta(t + \Delta T_{ij}) = \delta_{-\Delta T_{ij}}$ . Importantly,  $\Delta T_{ij}$  can then be directly related to a source azimuth  $\theta$  thanks to a relation of the form  $l_{ij}/c \cos \theta$ , with  $l_{ij}$  the interspace between the two considered microphones, when working in the farfield. If the two noise signals  $n_i(t)$  and  $n_j(t)$  are independent of  $s(t)$ , then the cross-correlation function  $R_{m_i m_j}$  comes as

$$R_{m_i m_j}(\tau) = E[m_i(t)m_j(t - \tau)] = (R_{ss} * h_{\mathbf{r}})(-\tau) + R_{n_i n_j}(\tau), \quad (17)$$

with  $R_{ss}$  the source autocorrelation function. Since  $h_{\mathbf{r}}(t) = \delta_{-\Delta T_{ij}}$ , and if the two signals  $n_i(t)$  and  $n_j(t)$  are independent, then one has

$$R_{m_i m_j}(\tau) = (R_{ss} * \delta_{-\Delta T_{ij}})(-\tau) = R_{ss}(\tau - \Delta T_{ij}), \quad (18)$$

bringing to the fore that  $R_{m_i m_j}$  is a temporally shifted version of  $R_{ss}$ . Since  $\forall \tau, R_{ss}(\tau) \leq R_{ss}(0)$ , then  $R_{m_i m_j}$  exhibits a maximum at  $\tau = \Delta T_{ij}$ . But in practice, as is the case for the MUSIC approach, the cross-correlation  $R_{m_i m_j}$  is not known since only one realization of the random signals  $m_i$  and  $m_j$  is available. The idea is then to build an estimation  $\hat{R}_{m_i m_j}$  of the cross-correlation, leading to the definition of the estimated TDOA

$$\widehat{\Delta T}_{ij} = \arg_{\tau} \max \left( \hat{R}_{m_i m_j}(\tau) \right). \quad (19)$$

Many cross-correlation estimators exist in the literature. One of the most used solution in Robotics consists in estimating the cross-correlation of filtered versions of the two signals  $m_i(t)$  and  $m_j(t)$ . This is obtained by introducing a function  $\Psi(f)$  weighting the frequency contributions of the two signals, the result being brought back in the time domain with an inverse Fourier transform, i.e.

$$\hat{R}_{m_i m_j}(\tau) = \int_{-\infty}^{+\infty} \Psi(f) \hat{S}_{m_i m_j}(f) e^{j2\pi f \tau} df, \quad (20)$$

where  $\hat{S}_{m_i m_j}(f)$  denotes the estimate of the cross-power spectral density function of the two signals  $m_i(t)$  and  $m_j(t)$ . Such estimators are known as *generalized cross-correlation* (GCC) techniques in the literature. Various different frequency weights have been proposed, most of them being listed and studied in [52]. Among them, one can cite the Roth [53], the Smoothed Coherent Transform (SCoT) [54], the Hannan-Thomson (HT) [55], or the Phase Transform (PhaT) processors. This last weighting is by far the most widely used in robotics, and is defined as

$$\Psi_{\text{PhaT}}(f) = \frac{1}{|\hat{S}_{m_i m_j}(f)|}. \quad (21)$$

From this definition,  $\hat{R}_{m_i m_j}(\tau)$  then comes as  $\hat{R}_{m_i m_j}(\tau) = \int_{-\infty}^{+\infty} e^{j\hat{\phi}(f)} e^{j2\pi f \tau} df$ , with  $\hat{\phi}(f)$  the phase of  $\hat{S}_{m_i m_j}(f)$ . In the ideal case when  $\hat{\phi}(f) \approx -2\pi f \Delta T_{ij}$ , then one gets

$\hat{R}_{m_i m_j}(\tau) \approx \delta(t - \Delta T_{ij})$ , i.e. the cross-correlation is different from zero only for  $\tau = \Delta T_{ij}$ , thus proving a very sharp estimation of the TDOA. Nevertheless, since the PhaT operation gives the same importance to all frequencies, it should not be used on narrowband signals, unless if some *a priori* on the frequency bandwidth of interest can be integrated. Such considerations will be discussed in §3.2.2.

All the aforementioned approaches mainly rely on a free field model, i.e. the  $i^{\text{th}}$  and  $j^{\text{th}}$  signals only differ in a delay  $\Delta T_{ij}$ . But as expected, the performances of this delay estimation highly degrade in the presence of reverberations, e.g. when working in a real robotic environment. This appears in the form of estimation outliers, which are all the more frequent as the reverberation time increases [56]. Additionally, the estimation strategy (19) leads in practice to a TDOA which is a multiple of the sampling frequency  $T_s$ , thus limiting the reachable angular resolution. For instance, for two microphones in the freefield spaced 16cm apart, and a sampling frequency  $f_s = 44.1\text{kHz}$ , this resolution spans from  $3^\circ$  (for a source facing the robot) to about  $18^\circ$  (for a sound at the left or right of the robot). Some interpolation strategies can nevertheless improve the resolution: interpolation with a parabola or an exponential function [57], or even interpolation of the whole cross-correlation function through *sinc* functions. Finally, as outlined for instance in (21), the correlation processor  $\Psi(f)$  must be estimated itself on the basis on an estimation of the cross-power spectral density function  $S_{m_i m_j}(f)$ . This can be achieved for instance by averaging short-term cross-periodograms (known as Welch's method [58]), for which the bias and variance have been studied with respect to the overlapping rate or the number of time window used for the estimation [59]. Similarly, it has been acknowledged that the duration of these windows has a critical effect on the accuracy of the TDOAs  $\Delta T_{ij}$  extracted from the cross-correlation peaks [60].

### 3.1.4. Beamforming based approaches

Among all the methods rooted in Signal Processing, those based on beamforming are probably the most used in Robotics. Their simplicity and low computational cost make them *a priori* well suited to this context. Yet, as will be shown, their performances strongly depend on the array characteristics, especially on its extent and number of microphones. This subsection then recalls some definitions and generalities on beamforming strategies used in robotics.

The term “beamforming” covers techniques to the combination of the signals coming from an array of discrete sensors, generally in order to focalize it to a specific direction of space  $\mathbf{r}_0$ . Typically, the signals  $m_n(t)$  spatially sampled at the  $N$  microphones locations  $n = 1, \dots, N$ , are processed by separate linear filters of impulse responses  $w_n(\mathbf{r}_0, t)$ . These filters are designed in such a way that the sum  $y_{\mathbf{r}_0}(t)$  of their outputs is the result of the spatial filtering described above. This principle is summarized in Fig. 6. On the basis on (8), the time relationship  $y_{\mathbf{r}_0}(t) = \sum_{n=1}^N w_n(\mathbf{r}_0, t) * m_n(t)$  can be turned into

$$Y_{\mathbf{r}_0}(k) = \sum_{s=1}^S D_{\mathbf{r}_0}(\mathbf{r}_s^s, k) S_s^0(k) + \sum_{n=1}^N W_n(\mathbf{r}_0, k) B_n(k), \quad (22)$$

with  $S_s^0(k)$  the frequency contribution at an arbitrary reference point 0 and due to the  $s^{\text{th}}$  source, with  $W_n(\mathbf{r}_0, k)$  the frequency response of the filter attached to the  $n^{\text{th}}$

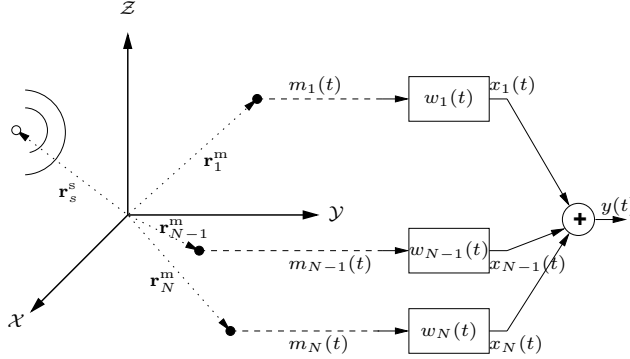


Figure 6: Basics of beamforming.

microphone, and

$$D_{\mathbf{r}_0}(\mathbf{r}, k) = \sum_{n=1}^N W_n(\mathbf{r}_0, k) V_n(\mathbf{r}, k). \quad (23)$$

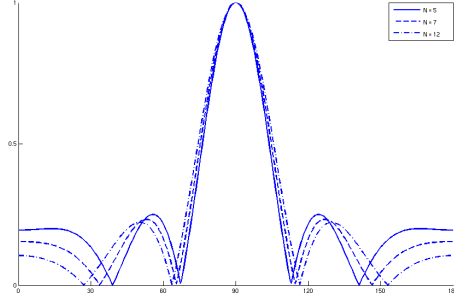
This last function of space and time variables is termed *array pattern*, or *beampattern*. It can be assimilated to a transfer function between the signal  $s_s^0(t)$  at the arbitrary reference position and caused by the  $s^{\text{th}}$  source emitting from position  $\mathbf{r}$ , to the beamformer output  $y(t)$ , and accounts for the amplification or attenuation of spatial areas. As (23) depends on  $V_n(\mathbf{r}, k)$ , this definition of the beampattern is valid both in the nearfield and in the farfield. Similarly to (11), the limiting expression  $D^\infty(\theta, \psi, k) = \lim_{r \rightarrow \infty} D_{\mathbf{r}_0}(\mathbf{r}, k)$  can be exhibited when the wavefronts are assumed planar. On this basis, an *energy map* of the environment  $E(\mathbf{r}, t)$  is then computed on a time window of length  $T$  along

$$E(\mathbf{r}, t) = \int_{t-T}^t |y_{\mathbf{r}}(\tau)|^2 d\tau, \quad (24)$$

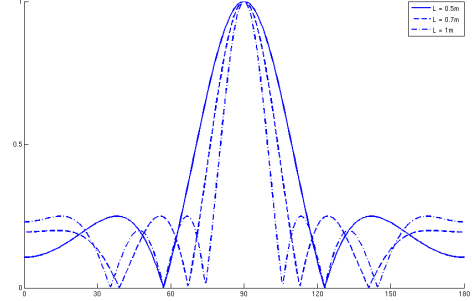
the sound sources positions being estimated by detecting the maximum of  $E(\mathbf{r}, t)$ . Practically, (24) is evaluated on a finite set of potential directions  $\mathbf{r}$ , see §3.2.2.

As already done in (11), and for the sake of simplicity, lets now consider a linear array, made up with  $N$  microphones aligned along the  $Z$ -axis and having the same interspace  $d$ , and whose abscissae  $z_n$  verify  $z_n = (n - \frac{N+1}{2})d$ . Consequently, the array length  $L = (N - 1)d$ . Such an array can be polarized towards a predefined azimuth  $\mathbf{r}_0 = \theta_0$  as soon as the filters  $w_n(\mathbf{r}_0, t)$  shown in Fig. 6 compensate the delays due to propagation so as to rephase the waves incoming from the DOA  $\mathbf{r}_0 = \theta_0$  prior to their summation. Under the planar wavefronts assumption, the transfer functions  $W_n(\mathbf{r}_0, k) = W_n(\theta_0, k)$  can be selected as  $W_n(\theta_0, k) = e^{jkz_n \cos \theta_0}$ , so that the farfield beampattern writes

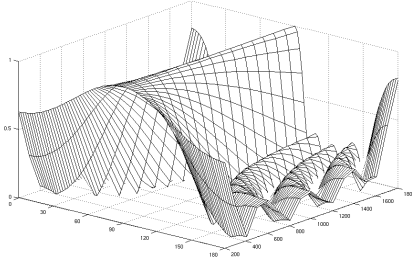
$$D_{\theta_0}^\infty(\theta) = \frac{\sin\left(\frac{\pi f}{c} N d (\cos \theta_0 - \cos \theta)\right)}{\sin\left(\frac{\pi f}{c} d (\cos \theta_0 - \cos \theta)\right)}. \quad (25)$$



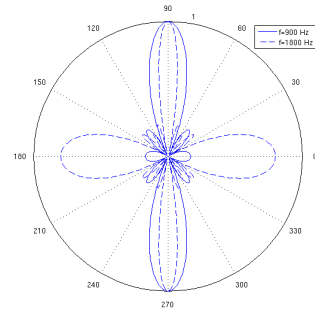
(a) Influence of the microphones number  $N$  of an array of fixed length  $L$  ( $f = 1kHz$ ,  $L = 0.7m$ ).



(b) Influence of the array length  $L$  for a fixed number  $N$  of microphones ( $f = 1kHz$ ,  $N = 5$ ).



(c) Normalized beam pattern as a function of  $\theta$  and  $f$  ( $N = 5$ ,  $L = 0.7m$ ).



(d) Illustration of the spatial aliasing.

Figure 7: Different beam patterns of a linear array. (a)&(b): Normalized beam patterns for various values of  $N$  and  $L$ . – (c)&(d): Influence of frequency of a beam pattern, for fixed  $N$  and  $L$ .

This so-called *conventional delay and sum beamforming* (DS-BF) strategy is by far the most used in robotics. For instance, Fig. 7(a) shows the module of (25) when considering  $\theta_0 = 90^\circ$ ,  $f = 1\text{ kHz}$  and  $L = 0.7\text{ m}$ . Several comments useful to Robotics can be deduced from this farfield array pattern expression in the following configurations:

1. variation of the microphones number  $N$ , for a fixed array length  $L$  and frequency  $k$  (or  $f$ );
2. change in the length  $L$ , for fixed  $N$  and  $k$ ;
3. modification of the frequency  $k$  for fixed  $N$  and  $L$ .

Such a study is fairly classical, see [61] [62], and is hereafter summarized for  $\theta_0 = 90^\circ$ . Increasing the number of microphones within a fixed-size array (scenario 1) leads to lower side lobes, see Fig. 7(a). The beam pattern corresponding to scenario 2 is shown on Fig. 7(b). The main lobe noticeably gets thinner as the array length increases. As a consequence, it may be necessary to mount a very large array on a robot in order to get a sharp focus towards a given direction of space. Embeddability constraints of

course prevent this, and thus limit the resolution of the whole acoustic sensor. Last, the third scenario is presented on Fig. 7(c). Keeping constant the microphones number and interspace, the main lobe width noticeably varies with the frequency  $f$ . The spatial resolution at low-frequency is poor, for high-wavelength waves are spatially oversampled by the array. A second phenomenon occurs at high frequencies: these are subject to aliasing, so that multiple replications of the main lobe appear. The spatial sampling of the wave must indeed obey a *Shannon spatial sampling theorem*, in that the maximal microphones interspace  $d$  must satisfy  $d < d_{max} = \frac{\lambda_{min}}{2} = \frac{c}{2f_{max}}$ , with  $f_{max}$  the maximal frequency in the wavefield. Fig. 7(d) illustrates the aliasing for an antenna made up with  $N = 5$  microphones spaced by  $d = 17.5$  cm, whose total length is then  $L = 0.7$  m.

This short overview of DS-BF performances demonstrate that being able to precisely focalize in a given direction requires a large array endowed with a lot of microphones, what may not be possible in a Robotics context. But even if it were, the resulting beampattern would be still a function of the source frequency, exhibiting a dramatical loss of resolution of the low frequencies. One solution could consist in ignoring such problematic frequency components, by filtering them out. But then it would be difficult to localize any speech signals, where most of the energy spreads from about 300Hz to 3.3kHz. Nevertheless, the computational cost of such approaches remains very low, with only  $N$  parallel filters running together, making them one of the most used localization techniques in Robotics.

### 3.2. Exploitation in Robotics

Now that the theoretical aspects have been overviewed, their applications to robotics are summarized. Following the lines of the above subsections, MUSIC, correlation and beamforming approaches are successively discussed.

#### 3.2.1. MUSIC

As shown in §3.1.2, MUSIC consists in computing –for only one frequency  $k_0$ – the so-called pseudo-spectrum  $h(r, \theta)$  defined in (15), from which the source position is extracted by isolating its maxima. Since the sources of interest in Robotics are mainly broadband, the approach needs to be extended to cope with multiple frequencies.

One of the first use of the MUSIC algorithm in robotics is [63]. Therein, an array of  $N = 8$  microphones, distributed on the periphery of the robot Jijo-2, enables the localization of vocal sources through an extension of the narrowband method to broadband signals. This extension, named SEVD-MUSIC (for Standard Eigen Value Decomposition-MUSIC), can be seen as “naive”, in that it closely follows the lines of the narrowband algorithm. First, the whole frequency range  $[k_L - k_H]$  of interest is partitioned into narrow frequency intervals, or “bins”, each one centered on  $k_b$ ,  $b = 1, \dots, B$ . The approximation of the covariance matrices  $\mathbf{\Gamma}_M(k_1), \dots, \mathbf{\Gamma}_M(k_B)$  are then computed, following a scheme similar to (14). From the subsequent eigendecomposition of each  $\widehat{\mathbf{\Gamma}}_M(k_b)$ , separate pseudo-spectra  $h_b(r, \theta)$  are determined,  $b = 1, \dots, B$ . The localization consists in isolating the maxima of the *average pseudo-spectrum*  $h_{Av}(\cdot)$

$$h_{Av}(r, \theta) = \frac{1}{B} \sum_{b=1}^B h_b(r, \theta). \quad (26)$$

Such a broadband extension is still in use in recent works [64, 65].

More recently, MUSIC received more attention from roboticists in order to deal with realistic scenarios possibly involving loud noise sources. In such a case, it can be difficult to easily identify the noise and signal spaces from the eigenvalue decomposition of the array cross-correlation matrix  $\mathbf{\Gamma}_M$ . For this reason, the GEVD-MUSIC (Generalized Eigen Value Decomposition-MUSIC) is proposed [66]. It consists in defining an additional freely-tunable correlation matrix  $\mathbf{\Gamma}_N$  for the frequency  $k_0$ , and solving the new GEVD problem

$$\mathbf{\Gamma}_M \mathbf{U}_n = \lambda_n \mathbf{\Gamma}_N \mathbf{U}_n, \quad (27)$$

where  $\mathbf{U}_n$  and  $\lambda_n$  depict the generalized eigenvectors and eigenvalues of the  $(\mathbf{\Gamma}_M, \mathbf{\Gamma}_N)$  matrix pencil respectively. The rest of the algorithm remains identical, since solving (27) allows to determine the noise and signal spaces, and then the computation of the pseudo-spectrum. Again, this operation is conducted along all frequency bins, to get an average pseudo-spectrum, along (26). The choice of the correlation matrix  $\mathbf{\Gamma}_N$  is free, but selecting  $\mathbf{\Gamma}_N = \mathbf{\Gamma}_B = E[\mathbf{B}\mathbf{B}^H]$  is a common choice which whitens the noise-related eigenvalues, and thus significantly eases the definition of the noise and signal spaces in the presence of loud noise sources. Interestingly, an extended, adaptive, version of GEVD-MUSIC has been proposed recently in [67]. Called iGEVD-MUSIC (for incremental GEVD-MUSIC), it consists in incrementally estimating the correlation matrix  $\mathbf{\Gamma}_N = E[\mathbf{B}\mathbf{B}^H]$  as a function of the current time frame. It then allows the use of the MUSIC algorithm in outdoor applications with drones, for which the level of the involved noises (ego-noise of the drone itself, and wind sound) is significative and especially dynamic [68]. But SEVD, GEVD or iGEVD approaches all suffer from the same problem: they are computationally expensive, with a high calculation cost for subspaces decomposition and for the pseudo-spectrum determination, both being generally performed on a frame-by-frame basis for real-time operations.

As a solution, the GSVD-MUSIC (Generalized Singular Value Decomposition-MUSIC) is proposed in [69]. As indicated by its name, it mainly relies on a generalized singular value decomposition, which consists in determining the left and right singular vectors  $\mathbf{U}_l$  and  $\mathbf{U}_r$  respectively, together with  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$  such that

$$\mathbf{\Gamma}_N^{-1} \mathbf{\Gamma}_M = \mathbf{U}_l \Lambda \mathbf{U}_r^H. \quad (28)$$

Once this decomposition is performed, the algorithm remains identical, with the left singular vectors and their corresponding singular values being used for the separation between the signal and noise spaces [69]. At the end, GSVD is shown to be computed almost 3 times quicker than GEVD, which is a critical improvement for real-time applications. But again, such a decomposition has to be performed for each frequency bin of interest. The contribution [70] exploits the idea of alignment as per [71], thus constituting a Coherent Broadband source localization algorithm (CB-MUSIC). Basically, the idea is to make the noise and signal spaces identical along all frequency bins through so-called *focalization matrices*  $T(\mathbf{r}, k_b)$  verifying  $T(\mathbf{r}, k_b) \mathcal{V}(\mathbf{r}, k) = \mathcal{V}(\mathbf{r}, k_0)$ , with  $k_0$  an arbitrary reference frequency. This way, the array vector at any frequency  $k$  is transformed into its value at frequency  $k_0$ . Thanks to this property, a unique correlation matrix gathering all the information along all frequency bins can be defined. Its generalized eigenvalue decomposition then allows the identification of the signal and noise spaces, and thus of the MUSIC pseudo-spectrum. In comparison with the other aforementioned approaches,

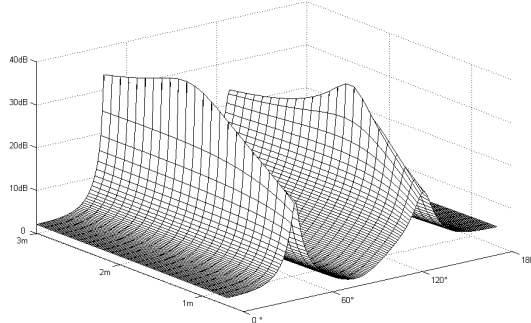


Figure 8: Typical MUSIC pseudo-spectrum, for two sources in the nearfield of a linear array.

only one generalized eigenvalue decomposition is necessary, thus limiting the computation cost of the method. Its implementation in a coherent beamspace paradigm is proposed along the lines of [72] and an original constructive method is proposed to the synthesis of focalization matrices, in a convex optimization setup. Besides, the approach is able to deal with reverberant robotics environments, since the statistical independence of the sources together with their mutual independence w.r.t. the noise can be relaxed.

All the approaches result in the computation of a pseudo-spectrum function. Such a function is depicted in Fig. 8, for a linear array and two independent sources placed respectively at  $(r, \theta) = (2m, 60^\circ)$  and  $(1m, 120^\circ)$ . As expected, the two very sharp peaks can be seen at the exact sources positions. But computing the pseudo-spectrum at each candidate source position can result in a high computational costs. A hierarchical strategy with a coarse-to-fine approach is proposed in [69] to solve this issue. Another hidden point concerns the source number, which must be known *before* identifying the noise and signal spaces, and thus before determining the broadband pseudo-spectrum. An information-theoretic approach, grounded on statistical identification —namely the Minimum Akaike Information Criterion Estimate (MAICE) defined in [73]— and relying on [71, 74] has been proposed in [75]. In addition to its sound theoretical bases, it has a very low computational cost and requires no prior threshold definition. Interestingly, the whole coherent beamspace MUSIC+MAICE detection and estimation has been implemented on a system-on-a-programmable-chip architecture [76].

### 3.2.2. Correlation-based approaches

*TDOA estimation.* The application to Robotics of correlation-based techniques presented in §3.1.3 is common since the beginning of Robot Audition. While the very first approaches were very naive, i.e. estimation of the TDOA by detecting the zero crossing points in the signals [36], the standard cross-correlation  $R_{m_i m_j}$  defined in (17) has been used in a lot of works. In [77], the intercorrelation is computed in order to infer the TDOAs between four microphones disposed on the vertices of a tetrahedron. The originality comes from the selection of the observation window: rather than computing the intercorrelation on the whole duration of the signals, a plain thresholding enables the detection of echoes-free temporal zones, onto which the TDOAs are determined. One can also cite [78], where a 4 microphones array is used to track a sounding docking station

outside the field of view of the robot. A slightly different application of the standard cross-correlation is proposed in [79], where the sound emitted by loudspeakers placed on the surface of a snake robot is used to estimate its posture using TDOA. Another traditional use of the standard cross-correlation consists in estimating the TDOA at the output of a filterbank [16]. Such an idea has been extensively used in a binaural context, where the filterbank is made of gammatone filters (see §2.2.2). This results in a TDOA function of the frequency, which is in essence analog to the IPD cue [34].

As already mentioned in 3.1.3, other strategies to cross-correlation computation exist. Among them, GCC techniques with the PhaT weighting function (GCC-PhaT) is by far the most used in a Robotics context. Its high temporal resolution in TDOA estimation justifies this choice, while it is known that this processor is highly sensitive to the length of the time windows used to estimate the cross-power spectral density function involved in (21). One can cite for instance [80], or [81] where the PhaT processor is exploited on a 24 evenly spaced microphones array fitted on the 3.2 m-long walls of a room which is visited by a tour-guide robot. More recent use of the PhaT approach can be cited: [82] where a triangular 3-microphone array is used to infer source location from short time observations so as to cope with the movement of the sound source or the robot; [83] presents an evaluation of various real-time sound localization approaches from a cubical 8-microphone array in which GCC-PhaT is compared with beamforming techniques; [84] proposes a robust approach to the acoustic perception of the presence of people from a pair of microphones. But GCC-PhaT only takes into account the phase of the perceived signals in the intercorrelation computation, giving the same importance to each frequency. As such a weighting does not differentiate the source and noise frequencies, the overall sensitivity of the method to noise is increased and voice localization becomes harder. As a solution, [85] defines an alternative processor which penalizes the frequencies at which the signal-to-noise ratio is low. This *Reliability-Weighted Phase Transform* (RWPhaT) strategy results in a new adaptive frequency weight  $\Psi(f)$ . This GCC strategy is still used in [86], on a 8-microphone array embedded on the Spartacus robot, to show the efficiency of a complete artificial audition system for speech recognition and dialogue management. Different adaptations of the PhaT processor have also been proposed. In [87], an eigenstructure-based GCC is outlined, based on the eigenvalue decomposition of the microphones auto-correlation matrix. Results show that the proposed processor exhibits less outliers than the traditional GCC-PhaT. In [88, 89], the GCC-PhaT- $\rho\gamma$  is proposed to deal with small SNR and large reverberation situations. Results demonstrate improvements w.r.t. the PhaT approach in terms of angular localization error, be the robot at rest or moving.

*From TDOA to localization.* Once the TDOAs have been computed by one of the above methods, the problem of localizing the source from their values must be addressed. For instance, consider a dipole in the farfield, made up with two microphones separated by a distance  $d_{ij}$ . In this planar wavefront case, the most direct approach to the determination of the azimuth  $\theta_s$  consists in inverting the formula  $\Delta T_{ij} = \frac{d_{ij}}{c} \cos \theta_s$ . This basic geometric rule is used in [90], the computed azimuths being involved into a neural network based sound source tracker. The same strategy is used in [78], or [84]. Following the same lines, one can deduce the cartesian coordinates  $\mathbf{r}^s = (u, v, w)$  of a source from the known positions  $\mathbf{r}_n^m = (x_n, y_n, z_n)$ ,  $n = 1, \dots, N$ , of the microphones constituting an array. If the propagation occurs in free space, the wavefronts impinging on the microphones are nested



spheres centered on the source. Under the assumption that each wavefront supports only one microphone and that the distance  $d$  from the source to the first receptor is related to the TDOAs  $\Delta T_{1n}$  between the 1<sup>st</sup> and every  $n^{\text{th}}$  microphone, the following holds:

$$\forall n \in [1, \dots, N], (x_n - u)^2 + (y_n - v)^2 + (z_n - w)^2 = (d + c\Delta T_{1n})^2. \quad (29)$$

After some manipulations, a matrix equation follows which leads to the unknowns  $(u, v, w, d)$ . This method is proposed in [91] to measure the time of flight of ultrasonic waves. Note that the antenna must hold at least 4 microphones in the planar case, 5 in the 3D case, otherwise the system is underdetermined. Unfortunately, the involved matrices may be ill-conditioned, so that very close TDOA values may lead to significantly different position estimates. This is why [92] proposes a simpler model, analogous to the one used in [77] and assuming planar waves. Noticing that the unit vector  $\nu = (u', v', w')$  pointing to the source—assumed to be at infinite distance—and the vector  $\mathbf{r}_{ij} = \mathbf{r}_j^m - \mathbf{r}_i^m$  connecting the  $i^{\text{th}}$  microphone to the  $j^{\text{th}}$  one satisfy

$$\forall n \in [1, \dots, N], \nu^T \mathbf{r}_{n1} = c\Delta T_{n1}, \quad (30)$$

$u', v', w'$  can be obtained through the resolution of a least square problem. In this approach, the matrix to be inverted depends solely on the microphones positions, and can therefore be tuned so as to improve its conditioning. Moreover, once the sensor geometry is fixed, the inverse matrix is constant and can thus be put in memory to reduce the necessary computations for localization. Nevertheless, the underlying propagation model assumes that planar wavefronts impinge on the antenna. As a solution, [87] recently proposed a generic extension of (30) to deal with the nearfield case, thus being able to estimate the distance to the source. Finally, a novel geometric formulation of the sound localization problem through TDOA is proposed in the very recent contribution [93], where an algebraic analysis and a global optimization solver are proposed for arbitrarily-shaped non-coplanar microphone arrays.

### 3.2.3. Beamforming

Among all the aforementioned strategies to sound source localization, beamforming remains probably the most exploited one in Robotics. As recalled in §3.1.4, beamformers are mainly designed to electronically polarize an array towards some specific DOA, and then to scan several directions of interest. An *acoustic energy map* can then be computed along (24), which is expected to be maximum at the actual sources DOAs. This strategy has been mainly coupled with Delay-And-Sum Beamformers (DS-BF) in a lot of contributions. Interestingly, the computational cost of DS-BF has been addressed in [94, 95] in two ways: first, the energy map is computed in the frequency domain through cross-correlations; next, the needed successive polarizations are performed towards directions defined by a recursive uniform icosahedron grid laid on a sphere. Other DOAs discretizations can be envisaged, depending upon the sensor shape and the number of test points, which lead to a tradeoff between the necessary computing power and the targeted resolution. But this conventional DS-BF strategy suffers from a lack of resolution in the polarization of low frequencies, together with the need of a high number of microphones, as demonstrated in §3.1.4. An example of such a DS-BF energy map for a short-length linear microphone array is shown in Fig. 9 (top-left) when trying to localize two speakers uttering from the azimuth 60° and 120°. Large main lobes regularly appear

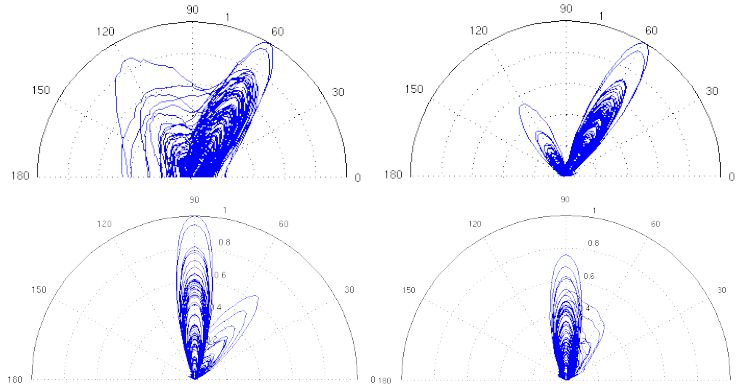


Figure 9: Acoustic energy maps of the environment (one curve per time snapshot) when using conventional beamformer (top left) or farfield frequency-invariant beamformer (top right and bottom left), see [96]. When used in the nearfield, a farfield frequency-invariant beamformer conducts to distorted energy maps (bottom right).

in the energy map, the two sound sources being then hardly spatially separable. Such a problem is often mentioned in the literature: in [97], where an array of 128 microphones spreaded into a room is used, the authors proposed to filter out all the frequencies below 500 Hz; the reference [98] gets close conclusions when simulating the 8-microphone antenna implemented on the small mobile platform EvBoy II: while the beampattern main lobe is thin enough for frequencies over 1 kHz, frequencies below 800 Hz cannot be exploited for localization; even with a three-ring 32-microphone array. [99] shows that the bad array directivity at low frequencies and the aliasing effect at low wavelengths conducts the localization to be performed only for frequencies between 1 and 2 kHz. More recent works still highlight this frequency limitation. For instance, [100] stated that if two sound sources are close to each other, false positive detections appear in the proposed system because of the wide directivity of DS-BF.

Different solutions have been proposed so far to deal with the low frequencies bad directivity of DS-BF approaches. In [100], an additional tracking step is introduced to reject the false detections. In [94], a probabilistic post-filtering of the acoustic energy map is performed, based on two simple short-term and mean-term estimators. Because of the temporal smoothing of the localization, a satisfactory robustness is achieved w.r.t. the actuators noise together with a reasonable computational complexity. An other solution consists in optimizing the array geometry so as to improve the consecutive beampattern. For instance, an evaluation index—relying on beampattern mainlobe width and sidelobes level measurements—is defined in [101, 102] so as to optimize the placement of 64 microphones over a 350-mm-diameter sphere. A valuable alternative may also consist in the synthesis of frequency-invariant broadband beamformers, as argued in [96]. Simulations of realistic scenarios entailing a 8-microphone linear array conclude to a significant improvement in the consequent acoustic maps, so that sources with close DOAs can be distinguished (see Fig. 9). Importantly, it is also established that the localization of sources emitting in the nearfield—e.g. at proximal human-robot interaction distance—is distorted if it entails a frequency-invariant beamformer designed under the

farfield assumption. An original nearfield frequency-invariant array pattern synthesis method is thus proposed, under the knowledge of the source range.

#### 4. Conclusion

Sound source localization methods developed in the Robotics community for the past 15 years have been reviewed in this paper. They can be partitioned into two classes. On the one hand, binaural techniques aim at reproducing artificially the human auditory system. The difficulty to exploit elementary acoustic cues has been underlined, together with the fundamental role of the head in the localization process. Though several propagation models have been proposed in the literature, the most basic of them are not sufficient to explain experimental measurements in an anechoic room. On the other hand, array processing techniques involving a larger number of microphones turn out to be intrinsically more accurate and robust. Different approaches were presented and their relevancy to Robotics was discussed. The extension of the high-resolution MUSIC method to broadband signals requires special care to cope with computational resources and the presence of noise in the environment. Correlation approaches to localization lead to accurate conclusions, yet they mainly assume planar wavefronts in order to limit the algorithmic complexity. Last, due to their versatility and low cost, beamforming based strategies are the most often used. However, the focalization of conventional beamformers is limited at low frequencies, so that alternative beamforming methods may be needed, and extreme care must be taken when dealing with nearfield sources. Importantly, a lot of the aforementioned contributions have been integrated within open software frameworks and hardware, making the most advanced approaches accessible to non-experts in the field. The most advanced solutions are the *HARK* (HRI-JP audition for robots with Kyoto University) software [103], the ManyEars framework [104], or the EAR (Embedded Audition for Robotics) system [105, 76].

Most of the cited works in this paper have only focused on the auditory scene analysis from a static view of the world. This idealized situation greatly eases the problem, while it is clear that speech and hearing takes place in a world where none of the static assumptions hold [106]. This is exactly what makes *Robot Audition* a Robotics problem on its own: the intrinsic mobility of modern robotics platform can be exploited to help in the environment analysis process. Actually, the Robotics Community has not extensively addressed this *active audition* topic, while it may constitute one of the most promising progress in embodied audition. Indeed, recent contributions in this field clearly demonstrate how the motion can be exploited together with the induced changes in the auditory perception to better the analysis, especially in the binaural framework. In this vein, [107] proposed a binaural sound localization system relying on the intersection of successive “cones of confusion” related to ITD measurements during a head movement. [108] also proposed to exploit the motion to dynamically reconfigure an array made of multiple microphones embedded on mobile robots to improve the sound localization. One can also mention [109] —who proposed a motion planning system whose objective is to maximize the effectiveness of a speech recognition module—, or [110, 111] —where the sound localization problem is rewritten in terms of a sensorimotor approach, with experiments made on the famous Psikharpax rat robot from the European FP7-ICT-IP ICEA (Integrating Cognition Emotion and Autonomy) project—. Stochastic filtering has also emerged as an ideal tool for sound localization and tracking during robot movement [112].

Recent contributions have proven the effectiveness of the approach in an active binaural experimental context, with the ability to cope with intermittent moving sources in the presence of false measurements [44, 43].

Providing robots with efficient and robust auditory functions will keep on being an exciting challenge during the forecoming years. Many problems which are considered as solved elsewhere have been renewed by the difficulties raised by the robotics context. The growing number of international projects dedicated to embodied audition clearly demonstrates the interest in this topic. Among them, one can cite the BINAAHR project (BINaural Active Audition for Humanoid Robots - French/Japan project funded by ANR and RSJ, ended in 2013), or the two just starting FP7 European projects EARS (Embodied Audition for RobotS) and TWO!EARS (Reading the world with TWO!EARS). All these forthcoming developments will be a source of stimulating discussions between the scientific communities of Acoustics, Signal Processing, Robotics, but also Physiology and Psychoacoustics. We then hope that this survey of existing approaches to the “low-level” stage of sound source localization will motivate new researchers to join the fertile field of Robot Audition.

## References

## References

- [1] S. Haykin, Z. Chen, The cocktail party problem, *Neural Comput.* 17 (9) (2005) 1875–1902.
- [2] S. Argentieri, A. Portello, M. Bernard, P. Danès, B. Gas, Binaural systems in robotics, in: J. Blauert (Ed.), *The Technology of Binaural Listening*, Springer, Berlin–Heidelberg–New York NY, 2013, Ch. 9, pp. 225–253.
- [3] T. Otsuka, K. Ishiguro, H. Sawada, H. Okuno, Unified auditory functions based on bayesian topic model, in: *Intelligent Robots and Systems (IROS)*, 2012 IEEE/RSJ International Conference on, 2012, pp. 2370–2376.
- [4] J. C. Middlebrooks, D. M. Green, Sound localization by human listeners, *Annual Reviews on Psychology* 42 (1991) 135–159.
- [5] M. Otani, S. Ise, Fast calculation system specialized for head-related transfer function based on boundary element method, *J Acoust Soc Am* 119 (5 Pt 1) (2006) 2589–98.
- [6] V. R. Algazi, R. O. Duda, D. M. Thompson, C. Avendano, The CIPIC HRTF database, in: *IEEE ASSP Workshop on Applications of Signal Processing to Audio and Acoustics*, 2001, pp. 99–102.
- [7] H. Wierstorf, M. Geier, S. Spors, A free database of head related impulse response measurements in the horizontal plane with multiple distances, in: *Audio Engineering Society Convention 130*, 2011.
- [8] C. I. Cheng, G. H. Wakefield, Introduction to head-related transfer functions (HRTFs): Representations of HRTFs in time, frequency, and space, *Journal of the Audio Engineering Society* 49 (4).
- [9] K. Nakadai, T. Lourens, H. Okuno, H. Kitano, Active audition for humanoid, in: *17th National Conference on Artificial Intelligence*, 2000, pp. 832–839.
- [10] K. Nakadai, H. Okuno, H. Kitano, Epipolar geometry based sound localization and extraction for humanoid audition, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, Vol. 3, 2001, pp. 1395–1401.
- [11] K. Nakadai, D. Matsuura, H. G. Okuno, H. Kitano, Applying scattering theory to robot audition system: Robust sound source localization and extraction, in: *Proceedings of the IEEE/RSJ Int. Conference on Intelligent Robots and Systems*, Vol. 2, 2003, pp. 1147–1152.
- [12] R. Duda, W. Martens, Range dependence of the response of a spherical head model, *Journal of the Acoustical Society of America* 104 (5) (1998) 3048–3058.
- [13] L. Rayleigh, On our perception of sound direction, *Philosophical magazine* 13 (74) (1907) 214–232.
- [14] R. A. Humanski, R. A. Butler, The contribution of the near and far ear toward localization of sound in the sagittal plane, *Journal of the Acoustical Society of America* 83 (10) (1988) 2300.

- [15] P. Zahorik, D. S. Brungart, A. W. Bronkhorst, Auditory distance perception in humans: A summary of past and present research, *Acta Acustica United With Acustica* 91 (3) (2005) 409–420.
- [16] K. Youssef, S. Argentieri, J.-L. Zarader, Towards a systematic study of binaural cues, in: *Intelligent Robots and Systems (IROS)*, 2012 IEEE/RSJ International Conference on, 2012, pp. 1004–1009.
- [17] P. I. M. Johannesma, The pre-response stimulus ensemble of neurons in the cochlear nucleus, in: *IPO (Ed.)*, *Symposium on Hearing Theory*, 1972, pp. 58–69.
- [18] R. D. Patterson, M. Allerhand, C. Giguère, Time domain modelling of peripheral auditory processing: a modular architecture and a software platform, *Journal of the Acoustical Society of America* 98 (1995) 1890–1894.
- [19] S. S. Stevens, J. Volkman, E. Newman, A scale for the measurement of the psychological magnitude pitch, *Journal of the Acoustical Society of America* 8 (3) (1937) 185–190.
- [20] P. Søndergaard, P. Majdak, The auditory modeling toolbox, in: J. Blauert (Ed.), *The Technology of Binaural Listening*, Springer, Berlin, Heidelberg, 2013, pp. 33–56.
- [21] R. Woodworth, H. Schlosberg, *Experimental psychology*, Holt, Rinehart and Winston NY (1962) 349–361.
- [22] A. A. Handzel, P. S. Krishnaprasad, Biomimetic sound-source localization, *IEEE Sensors Journal* 2 (2002) 607–616.
- [23] A. A. Handzel, S. B. Andersson, M. Gebremichael, P. Krishnaprasad, A biomimetic apparatus for sound-source localization, in: *IEEE Conference on Decision and Control*, Vol. 6, 2003, pp. 5879 – 5884.
- [24] K. Nakadai, H. G. Okuno, H. Kitano, Real-time sound source localization and separation for robot audition, in: *International Conference on Spoken Language Processing (ICSLP-2002)*, 2002, pp. 193–196.
- [25] R. Brooks, T. Senior, P. Uslenghi, The cog project: Building a humanoid robot, in: C. Nehaniv (Ed.), *Computations for Metaphors, Analogy, and Agents*, Springer Verlag, 1999, pp. 52–87.
- [26] J. Huang, N. Ohnishi, N. Sugie, Separation of Multiple Sound Sources by using Directional Information of Sound Source, Vol. 1, Springer Tokyo, 1997.
- [27] K. Nakadai, K. ichi Hidai, H. G. Okuno, H. Kitano, Real-time speaker localization and speech separation by audio-visual integration, in: *Proceedings of the IEEE International Conference on Robotics & Automation*, Vol. 1, 2002, pp. 1043–1049.
- [28] K. Nakadai, H. G. Okuno, H. Kitano, Auditory fovea based speech separation and its application to dialog system, in: *IEEE/RSJ International Conference on Intelligent Robots and System*, Vol. 2, 2002, pp. 1320–1325.
- [29] K. Youssef, S. Argentieri, J.-L. Zarader, A learning-based approach to robust binaural sound localization, in: *Intelligent Robots and Systems (IROS)*, 2013 IEEE/RSJ International Conference on, 2013, pp. 2927–2932.
- [30] E. Berglund, J. Sitte, Sound source localisation through active audition, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2005, pp. 509–514.
- [31] J. Hornstein, M. Lopes, J. Santos-Victor, F. Lacerda, Sound localization for humanoid robots - building audio-motor maps based on the hrtf, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2006, pp. 1170–1176.
- [32] H.-D. Kim, K. Komatani, T. Ogata, H. Okuno, Design and evaluation of two-channel-based sound source localization over entire azimuth range for moving talkers, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, Nice, France, 2008, pp. 2197–2203.
- [33] J. Liu, H. Erwin, S. Wermter, Mobile robot broadband sound localisation using a biologically inspired spiking neural network, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2008, pp. 2191–2196.
- [34] H. shik Kim, J. Choi, Binaural sound localization based on sparse coding and som, in: *Intelligent Robots and Systems*, 2009. *IROS 2009. IEEE/RSJ International Conference on*, 2009, pp. 2557–2562.
- [35] T. Rodemann, G. Ince, F. Joubin, C. Goerick, Using binaural and spectral cues for azimuth and elevation localization, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2008, pp. 2185–2190.
- [36] S. Cavaco, J. Hallam, A biologically plausible acoustic azimuth estimation system, in: *Third International Workshop on Computational Auditory Scene Analysis of the Sixteenth International Joint Conference on Artificial Intelligence*, 1999, pp. 78–87.
- [37] M. Kumon, T. Shimoda, R. Kohzawa, Z. Iwai, Audio servo for robotic systems with pinnae, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2005, pp. 885–890.
- [38] T. Shimoda, T. Nakashima, M. Kumon, R. Kohzawa, I. Mizumoto, Z. Iwai, Spectral cues for

- robust sound localization with pinnae, in: IEEE International Conference on Intelligent Robots and Systems, 2006, pp. 386–391.
- [39] A. Saxena, A. Ng, Learning sound location from a single microphone, in: IEEE International Conference on Robotics and Automation, 2009, pp. 1737–1742.
- [40] M. Kumon, Y. Noda, Active soft pinnae for robots, in: Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on, 2011, pp. 112–117.
- [41] J. Hebrank, D. Wright, Spectral cues used in the localization of sound sources on the median plane, *Journal of the Acoustical Society of America* 56 (6) (1974) 1829–1834.
- [42] E. Lopez-Poveda, R. Meddis, A physical model of sound diffraction and reflections in the human concha, *Journal of the Acoustical Society of America* 100 (5) (1996) 3248–3259.
- [43] I. Markovic, A. Portello, P. Danes, I. Petrovic, S. Argentieri, Active speaker localization with circular likelihoods and bootstrap filtering, in: Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on, 2013, pp. 2914–2920.
- [44] A. Portello, P. Danes, S. Argentieri, Active binaural localization of intermittent moving sources in the presence of false measurements, in: Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on, 2012, pp. 3294–3299.
- [45] T. Rodemann, A study on distance estimation in binaural sound localization, in: IEEE/RSJ International Conference on Intelligent Robots and Systems, 2010, pp. 425–430.
- [46] Y.-C. Lu, M. Cooke, Binaural estimation of sound source distance via the direct-to-reverberant energy ratio for static and moving sources, *IEEE Transactions on Audio, Speech, and Language Processing* 18 (7) (2010) 1793–1805.
- [47] S. Vesa, Binaural sound source distance learning in rooms, *IEEE Transactions on Audio, Speech, and Language Processing* 17 (8) (2009) 1498–1507.
- [48] H. L. Van Trees, *Optimum Array Processing, Vol. IV of Detection, Estimation, and Modulation Theory*, John Wiley & Sons, Inc., 2002.
- [49] R. Schmidt, Multiple emitter location and signal parameter estimation, *Antennas and Propagation, IEEE Transactions on* 34 (3) (1986) 276–280.
- [50] D. Purves, G. J. Augustine, D. Fitzpatrick, W. C. Hall, A.-S. LaMantia, J. O. McNamara, S. M. Williams, *Neuroscience, 3rd Edition*, Sinauer Associates, 2004.
- [51] L. A. Jeffress, A place theory of sound localization., *Journal of Comparative and Physiological Psychology* 41 (1948) 35–39.
- [52] C. Knapp, G. Carter, The generalized correlation method for estimation of time delay, *IEEE Transactions on Acoustics, Speech, and Signal Processing* 24 (4) (1976) 320–327.
- [53] P. R. Roth, Effective measurements using digital signal analysis, *Spectrum, IEEE* 8 (4) (1971) 62–70.
- [54] G. C. Carter, A. H. Nuttall, P. Cable, The smoothed coherence transform, *Proceedings of the IEEE* 61 (10) (1973) 1497–1498.
- [55] E. Hannan, P. Thomson, Estimating group delay, *Biometrika* 60 (2) (1973) 241 – 253.
- [56] T. Gustafsson, B. Rao, M. Trivedi, Source localization in reverberant environments: modeling and statistical analysis, *Speech and Audio Processing, IEEE Transactions on* 11 (6) (2003) 791–803.
- [57] T. May, S. van de Par, A. Kohlrausch, A probabilistic model for robust localization based on a binaural auditory front-end, *Audio, Speech, and Language Processing, IEEE Transactions on* 19 (1) (2011) 1–13.
- [58] P. Welch, The use of fast fourier transform for the estimation of power spectra: A method based on time averaging over short, modified periodograms, *Audio and Electroacoustics, IEEE Transactions on* 15 (2) (1967) 70 – 73.
- [59] G. Carter, C. Knapp, A. Nuttall, Estimation of the magnitude-squared coherence function via overlapped fast fourier transform processing, *Audio and Electroacoustics, IEEE Transactions on* 21 (4) (1973) 337 – 344.
- [60] M. Omologo, P. Svaizer, Acoustic event localization using a crosspower-spectrum phase based technique, in: IEEE International Conference on Acoustics, Speech, and Signal Processing, Vol. 2, 1994, pp. 273–276.
- [61] I. A. McCowan, Robust speech recognition using microphone arrays, Ph.D. thesis, Queensland University of Technology, Australia (2001).
- [62] S. Argentieri, P. Danès, P. Souères, Prototyping filter-sum beamformers for sound source localization in mobile robotics, in: IEEE International Conference on Robotics and Automation, 2005, pp. 3551–3556.
- [63] F. Asano, H. Asoh, T. Matsui, Sound source localization and signal separation for office robot jijo-2, in: IEEE/SICE/RSJ International Conference on Multisensor Fusion and Integration for

- Intelligent Systems, 1999, pp. 243–248.
- [64] C. Ishi, D. Liang, H. Ishiguro, N. Hagita, The effects of microphone array processing on pitch extraction in real noisy environments, in: Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on, 2011, pp. 550–555.
- [65] C. Ishi, J. Even, N. Hagita, Using multiple microphone arrays and reflections for 3d localization of sound sources, in: Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on, 2013, pp. 3937–3942.
- [66] K. Nakamura, K. Nakadai, F. Asano, G. Ince, Intelligent sound source localization and its application to multimodal human tracking, in: Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on, 2011, pp. 143–148.
- [67] K. Okutani, T. Yoshida, K. Nakamura, K. Nakadai, Outdoor auditory scene analysis using a moving microphone array embedded in a quadcopter, in: Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on, 2012, pp. 3288–3293.
- [68] K. Furukawa, K. Okutani, K. Nagira, T. Otsuka, K. Itoyama, K. Nakadai, H. Okuno, Noise correlation matrix estimation for improving sound source localization by multirotor uav, in: Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on, 2013, pp. 3943–3948.
- [69] K. Nakamura, K. Nakadai, G. Ince, Real-time super-resolution sound source localization for robots, in: Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on, 2012, pp. 694–699.
- [70] S. Argentieri, P. Danès, Broadband variations of the music high-resolution method for sound source localization in robotics, in: IEEE/RSJ International Conference on Intelligent Robots and Systems, 2007, pp. 2009–2014.
- [71] H. Wang, M. Kaveh, Coherent signal-subspace processing for the detection and estimation of angles of arrival of multiple wide-band sources, *IEEE Transactions on Acoustics, Speech, and Signal Processing* 33 (1985) 823–831.
- [72] D. B. Ward, T. D. Abhayapala, Range and bearing estimation of wideband sources using an orthogonal beamspace processing structure, in: IEEE International Conference on Acoustics, Speech, and Signal Processing, Vol. 2, 2004, pp. 109–112.
- [73] H. Akaike, A new look at the statistical model identification, *Automatic Control, IEEE Transactions on* 19 (6) (1974) 716–723.
- [74] M. Wax, T. Kailath, Detection of signals by information theoretic criteria, *Acoustics, Speech and Signal Processing, IEEE Transactions on* 33 (2) (1985) 387–392.
- [75] P. Danès, J. Bonnal, Information-theoretic detection of broadband sources in a coherent beamspace music scheme, in: Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on, 2010, pp. 1976–1981.
- [76] V. Lunati, J. Manhes, P. Danes, A versatile system-on-a-programmable-chip for array processing and binaural robot audition, in: Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on, 2012, pp. 998–1003.
- [77] F. Okuyama, J. ya Takayama, S. Ohyama, A. Kobayashi, A study on determination of a sound wave propagation direction for tracing a sound source, in: Proceedings of the 41st SICE Annual Conference, Vol. 2, 2002, pp. 1102–1104.
- [78] R. Luo, C. Huang, C. Huang, Search and track power charge docking station based on sound source for autonomous mobile robot applications, in: Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on, 2010, pp. 1347–1352.
- [79] Y. Bando, T. Mizumoto, K. Itoyama, K. Nakadai, H. Okuno, Posture estimation of hose-shaped robot using microphone array localization, in: Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on, 2013, pp. 3446–3451.
- [80] B. Mungamuru, P. Aarabi, Enhanced sound localization, *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on* 34 (3) (2004) 1526–1540.
- [81] Q. H. Wang, T. Ivanov, P. Aarabi, Acoustic robot navigation using distributed microphone arrays, *Information Fusion* 5 (2) (2004) 131–140.
- [82] C.-T. Kim, T.-Y. Choi, B. Choi, J.-J. Lee, Robust estimation of sound direction for robot interface, in: Robotics and Automation, 2008. ICRA 2008. IEEE International Conference on, 2008, pp. 3475–3480.
- [83] A. Badali, J.-M. Valin, F. Michaud, P. Aarabi, Evaluating real-time audio localization algorithms for artificial audition in robotics, in: Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on, 2009, pp. 2033–2038.
- [84] B. Hilsenbeck, N. Kirchner, Listening for people: Exploiting the spectral structure of speech to robustly perceive the presence of people, in: Intelligent Robots and Systems (IROS), 2011

- IEEE/RSJ International Conference on, 2011, pp. 2903–2909.
- [85] J.-M. Valin, F. Michaud, J. Rouat, Robust 3d localization and tracking of sound sources using beamforming and particle filtering, in: IEEE International Conference on Acoustics, Speech, and Signal Processing, 2006, pp. IV–IV.
- [86] M. Fréchette, D. Letourneau, J. Valin, F. Michaud, Integration of sound source localization and separation to improve dialogue management on a robot, in: Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on, 2012, pp. 2358–2363.
- [87] J.-S. Hu, C.-H. Yang, C.-K. Wang, Estimation of sound source number and directions under a multi-source environment, in: Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on, 2009, pp. 181–186.
- [88] H. Liu, M. Shen, Continuous sound source localization based on microphone array for mobile robots, in: Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on, 2010, pp. 4332–4339.
- [89] H. Liu, Z. Fu, X. Li, A two-layer probabilistic model based on time-delay compensation for binaural sound localization, in: Robotics and Automation (ICRA), 2013 IEEE International Conference on, 2013, pp. 2705–2712.
- [90] J. Murray, S. Wermter, H. Erwin, Auditory robotic tracking of sound sources using hybrid cross-correlation and recurrent networks, in: IEEE/RSJ International Conference on Intelligent Robots and Systems, 2005, pp. 891–896.
- [91] A. Mahajan, M. Walworth, 3-d position sensing using the differences in the time-of-flights from a wave source to various receivers, IEEE Transactions on Robotics and Automation 17 (1) (2001) 91–94.
- [92] J.-M. Valin, F. Michaud, J. Rouat, D. L’etouneau, Robust sound source localization using a microphone array on a mobile robot, in: Proceedings of the IEE/RSJ Int. Conference on Intelligent Robots and Systems, Vol. 2, 2003, pp. 1228–1233.
- [93] X. Alameda-Pineda, R. Horaud, A geometric approach to sound source localization from time-delay estimates, Audio, Speech, and Language Processing, IEEE/ACM Transactions on 22 (6) (2014) 1082–1095.
- [94] J.-M. Valin, F. Michaud, B. Hadjou, J. Rouat, Localization of simultaneous moving sound sources for mobile robot using a frequency-domain steered beamformer approach, in: Prodings of the IEEE International Conference on Robotics and Automation, Vol. 1, 2004, pp. 1033–1038.
- [95] J.-M. Valin, F. Michaud, J. Rouat, Robust localization and tracking of simultaneous moving sound sources using beamforming and particle filtering, Robotics and Autonomous Systems 55 (3) (2007) 216 – 228.
- [96] S. Argentieri, P. Danès, P. Souères, Modal analysis based beamforming for nearfield or farfield speaker localization in robotics, in: IEEE/RSJ International Conference on Intelligent Robots and Systems, 2006, pp. 866–871.
- [97] Y. Tamai, S. Kagami, H. Mizoguchi, Y. Amemiya, K. Nagashima, TachioTakano, Real-time 2 dimensional sound source localization by 128-channel huge microphone array, in: IEEE International Workshop on Robot and Human Interactive Communication, 2004, pp. 65–70.
- [98] L. Mattos, E. Grant, Passive sonar applications: Target tracking and navigation of an autonomous robot, in: IEEE International Conference on Robotics and Automation, Vol. 5, 2004, pp. 4265–4270.
- [99] Y. Tamai, Y. Sasaki, S. Kagami, H. Mizoguchi, Three ring microphone array for 3d sound localization and separation for mobile robot audition, in: IEEE/RSJ International Conference on Intelligent Robots and Systems, 2005, pp. 903–908.
- [100] Y. Sasaki, N. Hatao, K. Yoshii, S. Kagami, Nested igmm recognition and multiple hypothesis tracking of moving sound sources for mobile robot audition, in: Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on, 2013, pp. 3930–3936.
- [101] Y. Sasaki, T. Hujihara, S. Kagami, H. Mizoguchi, K. Oro, 32 Channel OmniDirectional Microphone Array Design and Implementation, Journal of Robotics and Mechatronics 23 (2011) 378–385.
- [102] Y. Sasaki, M. Kabasawa, S. Thompson, S. Kagami, K. Oro, Spherical microphone array for spatial sound localization for a mobile robot, in: Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on, 2012, pp. 713–718.
- [103] K. Nakadai, T. Takahashi, H. G. Okuno, H. Nakajima, Y. Hasegawa, H. Tsujino, Design and implementation of robot audition system HARK open source software for listening to three simultaneous speakers, Advanced Robotics 24 (5-6) (2010) 739–761. [arXiv:http://dx.doi.org/10.1163/016918610X493561](http://dx.doi.org/10.1163/016918610X493561),
- [104] F. Grondin, D. Ltourneau, F. Ferland, V. Rousseau, F. Michaud, The manyears open framework,



- Autonomous Robots 34 (3) (2013) 217–232.
- [105] J. Bonnal, S. Argentieri, P. Danès, J. Manhès, P. Souères, M. Renaud, The ear project, *Journal of the Robotics Society of Japan (RSJ)*, Special issue "Robot Audition" 28 (1) (2010) 10–13.
  - [106] M. Cooke, Y.-C. Lu, Y. Lu, R. P. Horaud, Active hearing, active speaking, in: *International Symposium on Auditory and Audiological Research*, Helsingor, Denmark, 2007, pp. 33–46.
  - [107] L. Kneip, C. Baumann, Binaural model for artificial spatial sound localization based on interaural time delays and movements of the interaural axis, *The Journal of the Acoustical Society of America* 124 (5) (2008) 3108–3119.
  - [108] E. Martinson, T. Apker, M. Bugajska, Optimizing a reconfigurable robotic microphone array, in: *Intelligent Robots and Systems (IROS)*, 2011 IEEE/RSJ International Conference on, 2011, pp. 125–130.
  - [109] M. Kumon, K. Fukushima, S. Kunimatsu, M. Ishitobi, Motion planning based on simultaneous perturbation stochastic approximation for mobile auditory robots, in: *Intelligent Robots and Systems (IROS)*, 2010 IEEE/RSJ International Conference on, 2010, pp. 431–436.
  - [110] M. Bernard, S. N’Guyen, P. Pirim, B. Gas, J.-A. Meyer, Phonotaxis behavior in the artificial rat *Psikharpax*, in: *Int. Symposium on Robotics and Intelligent Sensors (IRIS’2010)*, Nagoya, Japan, 2010, pp. 118–122.
  - [111] M. Bernard, P. Pirim, A. de Cheveigne, B. Gas, Sensorimotor learning of sound localization from an auditory evoked behavior, in: *Robotics and Automation (ICRA)*, 2012 IEEE International Conference on, 2012, pp. 91–96.
  - [112] Y.-C. Lu, M. Cooke, Motion strategies for binaural localisation of speech sources in azimuth and distance by artificial listeners, *Speech Communication*.