

Accelerated Speech Source Localization via a Hierarchical Search of Steered Response Power

Dmitry N. Zotkin*, Ramani Duraiswami

Perceptual Interfaces and Reality Laboratory, UMIACS
University of Maryland, College Park, MD 20742 USA

dz@umiacs.umd.edu, ramani@umiacs.umd.edu

Phone: (301)405-8753

Fax: (301)405-6707

Abstract

Accurate and fast localization of multiple speech sound sources is a problem of significant interest in applications such as conferencing systems. Recently approaches that are based on search for local peaks of the steered response power are becoming popular, despite their known computational expense. Based on the observation that the wavelengths of the sound from a speech source are comparable to the dimensions of the space being searched, and that the source is broadband, we develop an efficient search algorithm. Significant speedups are achieved by using coarse-to-fine strategies in both space and frequency. We present applications of the search algorithm to speed up simple delay-and-sum beamforming and SRP-PHAT source localization algorithms. A systematic series of comparisons with previous algorithms are made that show that the technique is much faster, is robust and accurate. The performance of the algorithm can further be improved by using constraints from computer vision.

Index Terms

EDICS Categories: 2-TRAN, 1-ENHA, 3-MCOM, 3-APPL.

I. INTRODUCTION

The inverse problem of localizing a source by using signal measurements at an array of sensors is an almost classical problem in signal processing. Along with the associated problem of beamforming, it has attracted the attention of many researchers. Our interest in this problem is in the context of localizing and possibly beamforming multiple sources of speech in a conferencing environment. As noted by Brandstein

This paper is an extended version of a conference paper [1]. Partial support of NSF awards 9987944 and 0086075 is gratefully acknowledged. We would also like to thank the anonymous referees for their comments and suggestions.

[2], many of the classical beamforming algorithms were motivated by applications in sonar or radar, rather than this particular problem, and consequently can perform poorly in the highly reverberant environments encountered in teleconferencing.

A recent book provides a very comprehensive introduction to the state-of-the-art in this field [3]. Generally, there are three classes of source localization algorithms: i) using steered beamformer energy response, ii) using high resolution spectral estimation, and iii) using time differences of arrival (TDOA) [4]. Some algorithms combine features from more than one class. We will focus on algorithms that fall in the first class. These algorithms, while capable of performing well, are usually slow since they involve a search for a peak of the response power. We show how the search can be performed efficiently in the case of speech sources in known rooms.

Inverse problems often exhibit ill-posed behavior in the sense of Hadamard [5], and their results are sensitive to noise in the data. In a reverberant environment, in addition, the data appears to be created by either the valid source position, or by any of a number of image sources induced by the scattering walls and surfaces. Thus, an additional element of ill-posedness is introduced in the problem, with the solution becoming multi-valued. Many algorithms are posed in the context of statistical signal processing, and do not treat this feature of the problem explicitly. A key to the solution of inverse problems is improved modeling that includes all available *a priori* information in the formulation. There has been some recent work in developing improved algorithms for this problem that use *a priori* information about the problem; for example, Brandstein presented algorithms for time delay estimation [2] and beamforming [6] that exploited knowledge about the pitch characteristics of speech. Our motivation here is similar, though we use different information.

A strategy that is often applied to resolve inverse problems is the iterative generation of a sequence of forward problems that might have created the data, with the best forward problem being taken to be one that minimizes an objective function. In the selection of the candidate forward problems, we can easily satisfy any *a priori* constraints. In the present case, the forward problem that is generating the data, is that there are speech source(s) in a room bounded by walls and other boundaries.

Application of this strategy to the present problem results essentially in the class of algorithm that we mentioned above – algorithms that steer the beamformer in various directions looking for the source and search for peaks in the output signal, also called SRP (steered response power) algorithms. The simplest (delay-and-sum) beamformer computes the propagation delays from the source position to each microphone and compensates for these delays in order to coherently sum the signals arising from the source position. More sophisticated beamformers filter the microphone signals in addition to delaying them. In any SRP algorithm the evaluation of an objective function has to be repeatedly performed. This

is usually the bottleneck in the performance of the algorithm.

As an illustration, one can build an energy map – the visual representation of variations of beamformer output energy versus the coordinates of the point that the beamformer is steered to (examples are shown later on). The source manifests itself as a peak in the energy map. The map depends on the array geometry and on the spectral content of the signal. The width of the peak in the energy map is generally smaller for higher-frequency sources. Increased microphone separation also decreases the width of the energy peak, which allows for higher localization accuracy and for better separation of sources (although increase of a microphone separation is limited by the appearance of spatial aliasing for high-frequency signals if the microphones are too far apart). Search for peaks in energy map is an obvious source localization algorithm. However, in real applications, the cost of computing the whole map at a fine enough resolution would be prohibitive, and some fast peak localization algorithm must be used. Traditional gradient descent [7] may be used if the search space is expected to contain one peak. If this is not true and the search space is multimodal, gradient descent is likely to find a local maximum. Then, many trials of gradient descent with random starts can be performed to improve the chances of finding the global maximum. Another fast search algorithm is stochastic region contraction, which is a general search algorithm for locating a global maximum of a multimodal function of many variables when the function satisfies certain conditions. It was successfully applied to microphone arrays, both in target localization [8] and in optimization of microphone placement [9], and allows significant reduction in number of objective function evaluations compared to repeated gradient descent with random starting points. Use of sequential Monte-Carlo methods [10], also known as particle filters, was also proposed for localization and tracking with microphone arrays (e.g., in [11], [12]). These algorithms are efficient since only a limited number of evaluations of the objective function are performed in the vicinity of the tracked position from previous frames.

This paper suggests a fast multi-level search strategy for an energy map using a coarse-to-fine paradigm in both the spatial and frequency domains, which we will call doubly hierarchical beamforming (DHBF). This strategy works because it uses the characteristics of the speech producing the objective function, e.g., the sound has characteristic wavelengths that are comparable to the dimension of the space that is being searched in this application (teleconferencing). The search algorithm we develop can be applied to any underlying energy function (we use two versions of SRP beamforming, simple delay-and-sum and phase-transform weighted, for illustration). It is not limited by the number of sources (though it does require them to be spatially separated to a certain degree and have similar power), or background noise structure, and has a predictable cost in terms of the number of evaluations of the objective function. It is particularly suitable for implementation in environments where there is prior knowledge of the spatial domain (e.g., the set of source locations being searched over can be restricted to an area bounded by room

walls, or the source is one of several objects detected via, say, computer vision, etc.). We show in the paper that the DHBF localization performance is comparable to other SRP based localization algorithms mentioned above and that the number of evaluations of objective function is significantly reduced.

II. BACKGROUND INFORMATION

We summarize here the *a priori* problem information and other background material known about the problem, and present some preliminary conclusions that can be drawn from them and be used to determine the coarse-to-fine strategy to speed up the steered response power search.

Spatial extent: The source occupies a region whose spatial extent is limited, so that in the search we can refine the search to a relatively small region. The environment is usually a workplace, a conference room or rarely an auditorium. In addition, sources are typically separated by distances of ~ 1 m, and will definitely be at least 0.3 m apart, since humans may be expected to be separated by this distance.

Nature of the speech signal: While the frequency range of human hearing extends from 20 Hz to 20 kHz, the sound produced by the human vocal tract has significantly less range, extending from about 100 Hz to 6 kHz. The spectral structure of the most energetic part of speech, the voiced phonemes (which include vowels and some consonants), consists mainly of a combination of integer multipliers of a fundamental frequency f_0 that lies between 80-200 Hz for males and 150-350 Hz for females. The voiced sounds constitute the low-frequency part of the speech spectrum. The other significant contribution are stop consonants and fricative consonants with their energy around 3-5 kHz. Overall, for a speech signal we can expect components in the range 100 Hz to 6 kHz [13].

f	λ	Feature	Remarks
20 Hz	17 m	Auditoriums	lower hearing limit
100 Hz	3.4 m	Conference rooms	speech lowest frequency
200 Hz	1.7 m	Rooms, human height	peak energy for speech
6 kHz	5.5 cm	Pinna Dimensions	speech highest frequency
20 kHz	1.7 cm	Concha size	upper hearing limit

Relationship between frequency and wavelength: The elementary equation $f\lambda = c$ indicates the relationship between frequency and wavelength. The wavelengths of audible sound are comparable to the dimensions of the space we live in and to our anatomical features. Humans use the spectral cues resulting from complex scattering of sound waves by objects with size comparable to the wavelength to determine the size of the environment and perform source localization [14]. Our goal is to exploit this relationship between audible speech frequency content and interesting spatial dimensions to develop fast search algorithms for locating sound sources. The table above presents some aspects of this relationship.

Delay and sum beamforming localization: Assume that the acoustic source that produces an acoustic signal $y(t)$ is located at point p and K receivers (microphones) are located at points q_1, \dots, q_K . The signal $s_m(t)$ at m^{th} microphone is given by

$$s_m(t) = y(t) \star h_m(q_m, p, t) + z_m(t), \quad (1)$$

where $h_m(q_m, p, t)$ is the room impulse response (RIR) function for the given positions of the source and the m^{th} microphone, star denotes convolution, and $z_m(t)$ is the combination of the channel noise and any environmental noise that is assumed to be independent at all microphones and not correlated with $y(t)$. As suggested in [4], we decompose the RIR into the direct arrival component (which simply consists of a single peak of amplitude r_m^{-1} at time $\tau_m = r_m/c$ where $r_m = \|p - q_m\|$ and c is the sound speed) and the rest of the RIR which we denote as $h_m^*(q_m, p, t)$, in which case equation (1) becomes

$$s_m(t) = r_m^{-1}y(t - \tau_m) + y(t) \star h_m^*(q_m, p, t) + z_m(t). \quad (2)$$

The received signals then explicitly contain delayed versions of the source signal plus its convolution with the rest of the RIR. We denote the time difference of arrival of a signal between receivers m and n as $\tau_{mn} = \tau_n - \tau_m$. The set of delays τ_{nm} can be associated with the location of the source. Another set of TDOAs $\hat{\tau}_{mn}$ is associated with the beamformer steering process. Given $\hat{\tau}_{mn}$, one can compute the output $s_B(t)$ of the delay-and-sum beamformer as

$$s_B(t) = \frac{1}{K} \sum_{m=1}^K s_m(t + \hat{\tau}_{mn}), \quad (3)$$

where n is a reference microphone which can be chosen to be the microphone closest to the position determined by the set of $\hat{\tau}_{mn}$ so that all $\hat{\tau}_{mn}$ are negative and the beamformer is causal. To steer the beamformer, we select $\hat{\tau}_{mn}$ corresponding to different positions in space, and if $\hat{\tau}_{mn} = \tau_{mn}$, the contributions from the source add coherently in the beamformed signal, resulting in unity gain for the source, whereas the signals of other (not steered to) sources and noise add incoherently, and their power will decrease on the average by a factor of K^{-1} . While this analysis is idealized since it assumes IID noise and neglects reverberation, it still provides useful insights into beamformer-based localization.

Equation (3) is in the time domain, but can also be expressed in the frequency domain. We recall that if a function $s(t)$ has Fourier transform $S(\omega)$, then time shifting of s by t_0 modifies its Fourier transform as $s(t - t_0) \Leftrightarrow S(\omega)e^{-j\omega t_0}$, where $j = \sqrt{-1}$. Equation (3) then becomes

$$S_B(\omega) = \frac{1}{K} \sum_{m=1}^K S_m(\omega)e^{j\omega\hat{\tau}_{mn}}, \quad (4)$$

and the output power up to a scale factor can be expressed as

$$P_B(\omega) = \sum_{m=1}^K \sum_{n=1}^K S_m(\omega)S_n^*(\omega)e^{j\omega\hat{\tau}_{mn}}. \quad (5)$$

A simple localization strategy can then be suggested by searching through the space of $\hat{\tau}_{mn}$ for an energy peak in the output $P_B(\omega)$. The search is usually performed in some organized fashion through the possible locations of sound source in 3D space, generating sets of TDOA corresponding to the spatial locations by trivial geometric computations. We describe below existing ways to improve robustness and speed of beamformer-based localization and suggest a novel fast search strategy based on hierarchical subdivision of space and frequency.

SRP-PHAT localization: In localization algorithms that use purely TDOAs to localize the sound source, the TDOAs are usually obtained using a generalized cross-correlation (GCC) between signals s_m and s_n acquired at the m^{th} and n^{th} sensors respectively [15]. Denote by $r_{mn}(\tau)$ the GCC of $s_n(t)$ and $s_m(t)$ and its Fourier transform by $R_{mn}(\omega)$. Then,

$$R_{mn}(\omega) = W_{mn}(\omega)S_m(\omega)S_n^*(\omega), \quad (6)$$

where $W_{mn}(\omega)$ is a weighting function. Ideally, $r_{mn}(\tau)$ (computed as the inverse Fourier transform of $R_{mn}(\omega)$) will have a peak at the true TDOA between sensors m and n (τ_{mn}). In practice, many factors such as noise, finite sampling rate, interfering sources and reverberation might affect the position of the peak. The phase transform (PHAT) weighting function was introduced in [15]:

$$W_{mn}(\omega) = |S_m(\omega)S_n^*(\omega)|^{-1}. \quad (7)$$

The PHAT weighting places equal importance on each frequency by dividing the spectrum by its magnitude. It was later shown [4], [16], [17], [18] that it is more robust and reliable in realistic reverberant conditions than other weighting functions designed to be statistically optimal under specific non-reverberant noise conditions. The SRP-PHAT algorithm [4], [19] applies the PHAT weighting in the context of the filter-and-sum beamformer. The power in its output $P_B(\omega)$ is given by

$$P_B(\omega) = \sum_{m=1}^K \sum_{n=1}^K W_{mn}(\omega)S_m(\omega)S_n^*(\omega)e^{j\omega\hat{\tau}_{mn}}, \quad (8)$$

which generalizes equation (5) by including the non-negative weighting function. In SRP-PHAT the function (7) is used to increase robustness in reverberant environments. As with simple delay-and-sum beamforming, the total power in the output of the beamformer $P_B = \int P_B(\omega)d\omega$ depends on how the beamformer is steered. The steering to a position p is done by computing the set of TDOA τ_{mn} that would have been observed for a source at that position from the known geometry of the microphone array. While acoustical sources in space produce peaks in the beamformer's output, the search for these peaks can be complicated by multiple peaks in the search space in case of SRP-PHAT algorithm. A sophisticated search mechanism is required to ensure successful localization of the global maximum.

Time delay imprecision: Given four or more receivers every point in physical space (x, y, z) can be mapped to a point in delay space $(\tau_{12}, \tau_{13}, \tau_{14}, \dots, \tau_{1K})$. The inverse map (from delays to source

location), is non-linear and ill-posed as discussed earlier. Here we consider the effect of errors in time-delays caused by an incorrect hypothesis of the source location on the computed steered response power. If the error in phase is small, then the coherence in the signals being added won't be completely destroyed, though incoherent components will also get an increase in their energy.

Consider the simplest beamformer consisting of two microphones separated by distance $2d$. If the beamformer main lobe is aligned perfectly on the source, the phase difference between signals at two microphones after appropriate time delays is zero and the beamformer gain is equal to 1 for the source whereas the gain for the incoherent noise can be expected to be equal to $\sqrt{2}/2$ and the power of the incoherent noise is halved. Simulation shows that as the misalignment of the beamformer (and consequently the phase difference between two signals) increases, the output power of the source will decrease and when the phase difference is equal to approximately $2\pi/5$ the power of the source drops close to the expected power of incoherent noise. Thus, it can be reasonably expected that when this simplest beamformer is misaligned by the distance which results in phase difference of less than $2\pi/5$ between channels, there still will be gain for the desired source compared to the incoherent part. It can also be shown that if a source located at a distance l from the array of diameter d is shifted by a distance b (which is the same as having the array misaligned by the same distance), the phase difference introduced by the shift is no more than $4\pi db/l\lambda$ (the phase shift is maximal when the source is on the axis of the array and movement is performed parallel to the array plane), where λ is the wavelength of the source. Thus, it is conservative in case of far-field source (when $l \gg d$) to estimate that an error in the source position of $\lambda/5$ will still result in a coherent gain in the beamformed signal. This is confirmed by tests with actual speech signals to our array. We refer to this result as the *imprecision heuristic*.

Another way to look at this result is to plot the spatial width of the beamformed signal peak as a function of the source frequency. From simulations performed by mixing actual room recordings of speech we see that there is an inverse relationship between the peak width in the energy map and sound wavelength. The peak in the energy map has an FWHM (full width at half maximum) of approximately $2\lambda/5$ for our array configuration, consistent with the heuristic. For the frequency of 150-160 Hz ($\sim f_0$) we get $2\lambda/5 \approx 0.8m$.

III. HIERARCHICAL LOCALIZATION ALGORITHM

Let there be K microphones and let us consider beamforming a data frame sampled at a frequency higher than the Nyquist rate for the highest frequency in the signal. We restrict our consideration first to the two-dimensional case and search for the direction of arrivals information, ignoring range for now. We divide the space of search parameters (azimuth φ and elevation θ) into quadrants of dimension corresponding to the coarsest size that is consistent with a gain in energy at the lowest frequency. We

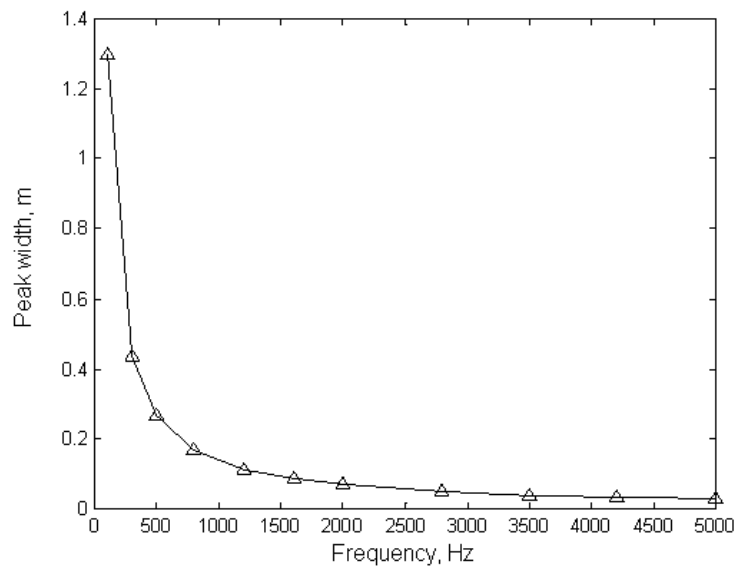


Fig. 1. Beamformer peak width as a function of frequency.

then evaluate the objective function (energy in the beamformed signal) at the centers of these quadrants, creating a coarse-level energy map, and search for local maxima in it. Each local maximum is tagged for further consideration. All tagged quadrants are further subdivided in the next pass of the algorithm.

There are two issues that must be fixed with this approach. First, at the coarse level we must restrict the beamforming to frequencies that are guaranteed to see an improvement in their power despite inaccurate steering (to the quadrant center instead of the true source position). In addition, performing beamforming using either (4) or (8) on the full signals for even the relatively few coarse quadrants will be uneconomical. A first possibility to fix this problem is decimation of the signal in the time domain. However, one quickly realizes that this leads to significant aliasing. An approach which achieves both goals is lowpass beamforming in the frequency domain which can be done quickly [20] simply by computing beamforming power only in those frequency bins where it is necessary.

In this approach we compute the FFTs of each of the received signals at the K channels. We decimate the signal in the frequency domain (performing a lowpass operation) with a cutoff frequency determined by the quadrant size at the coarsest level and compute the power at the centers of the quadrants. Let there be k sources, leading to k tagged quadrants. These are further recursively divided, using the quadtree data structure described below, and the power is computed at child nodes with new cutoff frequency twice as high at each step (because the length of the quadrant side gets divided by two at each level of quadtree subdivision). The child node with the highest energy level is selected as the most probable source position, and the subdivision is repeated recursively until desired node level and desired precision is achieved.

A. *Quadrees and octrees*

To perform a hierarchical division of space we use quadrees. A quadtree is a data structure for hierarchical representation and processing of spatial data in 2D, typically organized as follows. The root of a quadtree is associated with a 2D region bounded by a parallelogram (often a square). This parent region is subdivided into four similar equally sized quadrants, each carrying more specific information about its portion of space. Each child quadrant is, in turn, recursively partitioned into four children. The process of subdivision can continue infinitely, but in practice it has to stop (e.g., when further subdivisions will not significantly help the search). The hierarchical nature of quadrees allows one to efficiently represent and search data distributed in 2D space, which is possible because the size of a quadtree representing a 2D region is $O(s)$ where s is the perimeter of the object. Algorithms that execute on quadrees rather than on pixel arrays have running times proportional to the number of blocks in the quadtree.

Similarly, the octree is the data structure representing 3D volumetric data, with eight children per node and number of nodes proportional to the area of the object's surface. Use of quadrees and octrees leads to the dimension reduction effect, i.e. an octree algorithm applied to a 3D problem is analogous to an array-based algorithm in 2D [21].

B. *Implementation*

The algorithm for localization of multiple sound sources using the proposed doubly-hierarchical search on steered response power function proceeds by first dividing the search area into quadrants forming a coarse $L \times L$ search grid. When it is known that the source(s) will be located in front of the microphone array (for example, when the array is on a wall) the search area can be defined as a square in the DOA space $\varphi \in [-\pi/2, \pi/2], \theta \in [-\pi/2, \pi/2]$. The quadrant size is chosen sufficiently small so that it is unlikely that two sound sources share the same quadrant and that local maxima of the energy can be computed. The latter condition is ensured by the $\lambda/5$ heuristic. Note that the quadrants in DOA space are not rectangular in the world coordinate frame, and the cutoff frequency and the initial grid size have to be selected by re-mapping DOA grid onto a (x, y, z) grid using the maximum possible source depth (determined by known room dimensions). The nodes that correspond to the local maxima in the constructed energy map are selected for further processing. Every node is recursively searched by partitioning it into four children, and the child with the maximum energy level is selected for subdivision at the next level. The recursion terminates when the quadtree branch reaches a certain depth, corresponding to a fixed minimal quadrant size. The procedure is repeated for about 10 levels which yields a quadrant of size of about 1 cm (when recomputed to spatial units) in our implementation. The center of the maximal energy quadrant at the deepest level is output as the source position.

A problem that may arise is that the actual peak may lie at the boundary of a quadrant in DOA space

resulting in possible mislocalization of the energy maximum during the coarsest stage of algorithm. To avoid this, we perform a simple check at the last step. If at the end of the recursive search within a quadrant the peak is localized at the boundary P of the original coarse quadrant V , a search is also performed in the neighboring coarse quadrant V' that shares P with V . In practice this rarely happens.

The algorithm as described above is developed for 2D sound source localization. Sometimes it is desirable to perform full 3D localization with one or more arrays [22], and our search algorithm can be adapted for this situation. In this case, the search algorithm is executed in a 3D space using an octree based peak search instead of a quadtree, which directly determines the Cartesian coordinates of the sound source(s) relative to the arrays.

IV. PERFORMANCE EVALUATION OF THE ALGORITHM

We compare the developed search algorithm with other steered response power algorithms. The algorithms we test are full search [4], repeated gradient descent [7], stochastic region contraction [9] and DHBF. We will define them with one-letter designations (for the plots) as follows:

- (F) Full search for the SRP-PHAT energy maximum over the 1024-by-1024 grid of all possible DOAs in the $Z > 0$ hemisphere;
- (G) Repeated gradient descent with $N = 1000$ different starting points (we found this number is the minimum sufficient for repeatable and robust localization of the global maximum);
- (S) Stochastic region contraction with the parameters suggested in [9] ($K_1 = 100$, $K_2 = 400$, $\hat{K} = 8$);
- (D) The proposed DHBF algorithm with 16-by-16 initial grid size and 1 kHz initial high frequency cutoff, searching up to the level 10 of quadtree subdivision.

A. Multisource energy map illustration

Let us start by observing some interesting properties of the energy map derived using the SRP-PHAT energy function. In Figure 2, energy in the beamformed signal is plotted as a pixel intensity on a two-dimensional plane with the horizontal axis being the azimuth $\varphi \in [-\pi/2, \pi/2]$ and the vertical axis the elevation $\theta \in [-\pi/2, \pi/2]$ (so essentially a hemisphere of directions with positive Z is plotted). From here on we will use simulated results on a circular planar array of 7 microphones with one microphone at the center and the other 6 located at the circumference of the array to match a real array we have. The radius of the array is 0.3 m. Assume the origin of the coordinate system at the center of the array with the X and Y axis in the plane of array and Z axis orthogonal to it, and the point with zero azimuth and zero elevation being directly in front of the array on the positive Z axis. The sources in these plots are placed at $(-3.3, 1.7, 2.5)$, $(2.5, -2.7, 3.5)$, and $(0.0, 4.2, 3.0)$ meters.

In Figure 2a, the sources are well-defined and sharp. However, the number of local maxima is much more than the number of sources. Each bright curved line is formed by a projection of the TDOA locus

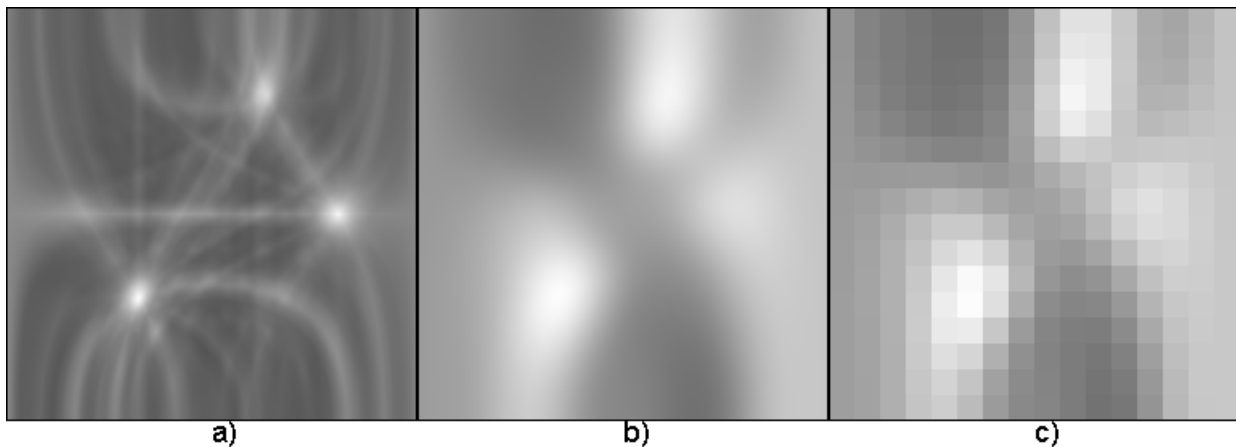


Fig. 2. a) Sample energy map (energy vs azimuth and elevation) obtained using SRP-PHAT algorithm with three sources, resolution 512x512. b) Same map for the signals lowpassed at 2 kHz. c) Same map for the signals lowpassed at 2 kHz, resolution 16x16.

(a hyperboloid for a given microphone pair) into the (φ, θ) space. This projection can be viewed as taking an intersection of the equi-TDOA surface centered at the point between the microphones in the pair with the sphere of some radius ρ centered at the center of the array and letting $\rho \rightarrow \infty$. As $\rho \rightarrow \infty$, the intersection locus describes all directions of arrivals that satisfy the observed TDOA. Lines corresponding to the same source intersect at the same spot. This is the visual illustration to the source localization by intersection of cones defined by each TDOA.

It is difficult to localize multiple sources in the map because of the many local maxima. In contrast, in Figure 2b the same map is plotted with the signal lowpassed at 2 kHz. All fine details are removed, and the map consists of three broad peaks located roughly at the source positions. This can be expected because of the relationship between peak width and frequency – the lower the frequency, the higher the spatial shift that will keep the source still in partial focus, and fine details are not present in the map because high frequencies that contribute to fast spatial variations of the beamformer output due to small wavelengths are missing. The rationale behind the coarse map of lowpassed signals as an initial step of our localization algorithm is that if the map in Figure 2b is downsampled (Figure 2c), the resulting coarse map has peaks in approximately correct positions, which are used as starting points for a recursive search.

Bandlimiting at low resolution is thus crucial to avoid misleading peaks and obtain correct initial estimates of source locations. It is also important to keep the signal bandlimited to appropriate frequencies determined by the quadrant size while performing a hierarchical subdivision of space. Further this strategy uses the important prior information that the peaks that arise at each frequency in a spatial neighborhood are caused by the same physical source.

B. Accuracy in reverberant environments

First, we measured and compared performance of several algorithms under simulated reverberation which is perhaps more degrading to the performance of localization than noise alone. We used clean speech (utterance of ten consecutive digits “zero, one, two, three...” with short pauses between the digits) recorded by a microphone placed close to the speaker’s mouth as a source signal and simulated reverberation using a simple image model [23] in a rectangular room. In the model, a regular lattice of virtual sources is created, which represents the reflections of the acoustic source in room walls (including floor and ceiling). A room transfer function (alternatively, room impulse response, or RIR) function can be computed using the image model [24]. Microphone outputs are computed by convolving the source waveform with the appropriate RIRs for each microphone position. The sampling frequency was set to 40.44 kHz to match the setup used in the real experiments (described later). The simulated room has dimensions $5 \times 2.5 \times 4$ m. The origin of the coordinate system is placed at the center of the room, and the Y axis is vertical so that the coordinates of one of the room corners is $(2.5, 1.25, 2.0)$. The center of the microphone array is at $(0.0, 0.0, -2.0)$ and the array is a 0.3 m radius array with 6 microphones equispaced on the array circumference and one microphone in the center. The RIR was computed up to the 24th reflection and lowpassed with cutoff frequency of 100 Hz as described in [23]. We ran several simulations with different wall reflective properties to simulate different reverberation times. Within each simulation, we placed the source at a random point within the room subject to constraints $X \in [-2.0, 2.0]$, $Y \in [-1.0, 1.0]$, $Z \in [-1.0, 1.0]$ and we tested with 64 points for each reverberation condition. We used non-overlapping 50 ms frames, and from a 5 s utterance 33 non-silence frames were selected for processing.

From the known source position, the correct DOA was computed and compared to the DOA estimate produced by the four search algorithms. The beamformed energy in all algorithms was computed using the SRP-PHAT technique, with frequency components from 300 Hz to 11 kHz included in the computations, except for the DHBF algorithm, where the high cutoff is adjusted during the search as described previously. We present experimental data in concise form here. In the left part of Figure 3, we show the average localization error for the four algorithms for 5 different reverberation times (90, 120, 150, 180, and 210 milliseconds, respectively) in the top histogram. Percentages of successful localizations are shown in the bottom histogram (where “successful localization” is defined as localization within 10 degrees of the true DOA). On the right, we show the same plots with all frames where full search fails to localize the source omitted (see discussion below), in which case full search obviously shows 100% correct localization.

To our surprise, in initial simulations we found that the DHBF, in addition to its speed, was superior to all other methods in terms of average localization error, even to the full search! This means that under simulated reverberation there exist some cases when the maximum peak of energy does not coincide with

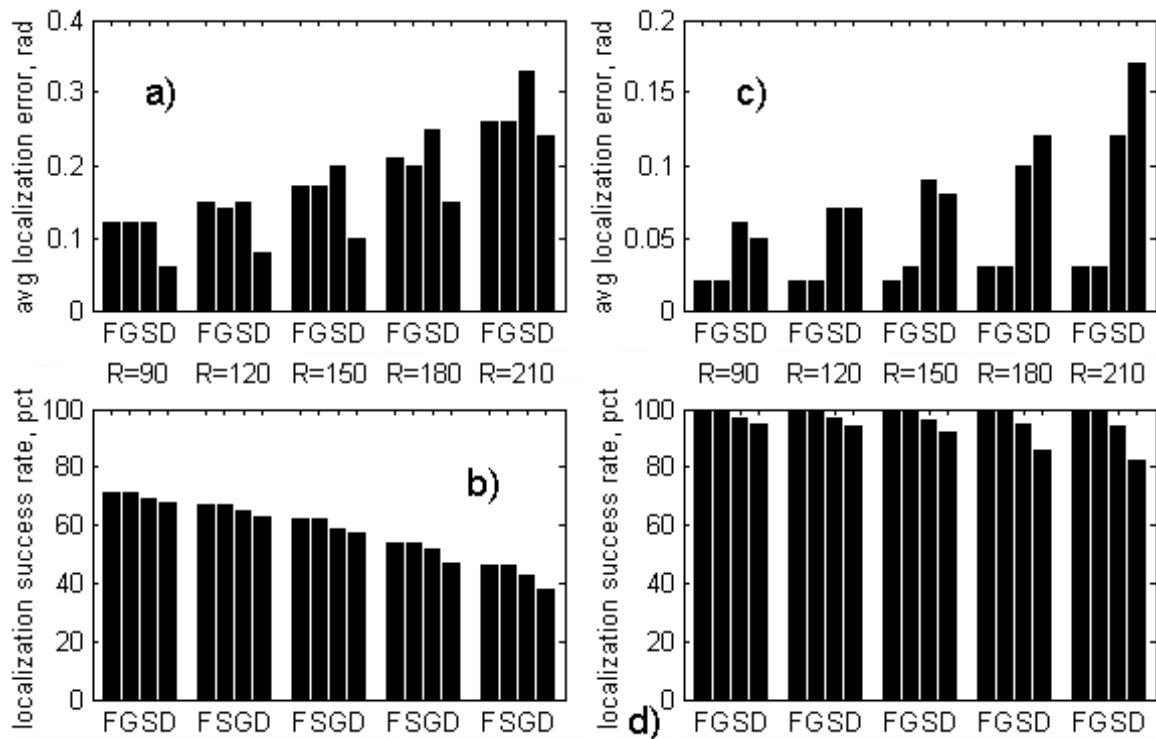


Fig. 3. a) Average localization error for 4 algorithms. b) Percentage of correct localization for 4 algorithms. c) Average localization error with data frames where full search fails omitted. d) Percentage of correct localization under the same condition. R is the reverberation time in milliseconds.

the true DOA (and thus full search, which presumably should be the most robust method, fails). This happens in about 30% of all cases with the smallest reverberation time (90 ms) and in about 55% of all cases with the largest one (210 ms). Examination of the energy maps in such cases shows that generally the highest (false) peak is located close to the $\varphi = \theta = 0$ DOA, in which case all search methods except DHBF are distracted by it. On the other hand, the lowpassed energy map that is the starting point for DHBF does not contain the false peak, and thus the DHBF search is initialized approximately correctly on the source and stays in the vicinity of the initialization point during the refinement stage. This results in a DOA estimation that is close to the correct DOA.

The appearance of a false peak located around the center of room is due to several factors, in particular to the symmetric arrangement of the microphones in the array and to the effect of the windowing operation, which causes a bias in Fourier transform coefficient phases (moreover, remember that SRP-PHAT operates essentially only on the phases of the coefficients and thus is particularly sensitive to the bias). It is also enhanced by the idealized room reverberation model with a regular lattice of virtual source. In real conditions, slight asymmetries of the array and room and objects present in the room are likely to diminish this effect. Indeed, in real experiments we did observe some data frames with such behavior (in which

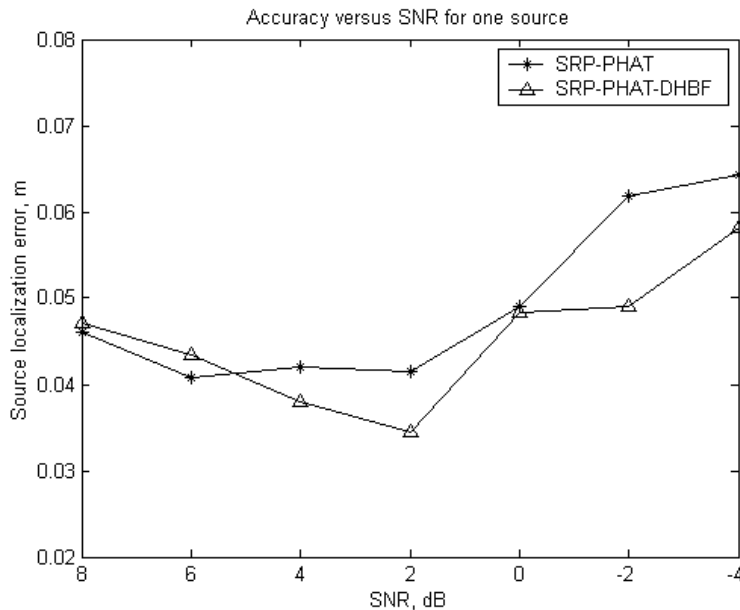


Fig. 4. Localization error versus SNR for SRP-PHAT and SRP-PHAT-DHBF.

all methods but DHBF find a false peak close to the center of room) but in reduced proportion compared to the simulated cases. We also performed preliminary simulated experiments using randomly scattered microphone positions in the array and smoother windowing operation, and both of these somewhat decrease the frequency of the false peak appearance but do not completely eliminate such cases. This topic is a subject of our future research.

In summary, we see that the localization performance of DHBF is comparable to the performance of other search algorithms in case of simulated reverberation. As the reverberation time increases, the DHBF performance somewhat worsens if the energy map has a pronounced peak at the true DOA. However, this is partially offset by good DHBF performance if there is a false energy peak near the center of the room. Moreover, the algorithm that is closest to DHBF in terms of speed, stochastic region contraction, is also worse in performance compared with the repeated gradient descent and full search.

C. Accuracy in noisy environments

We performed testing of the localization accuracy of the SRP-PHAT algorithm with and without our search algorithm under different noise conditions. We synthesized an acoustic scene in which one source was located at a known position and contaminated the signal in each channel with white noise to achieve different SNRs. In Figure 4, we show the average localization error in 20 trials with random positioning of the source for the original SRP-PHAT algorithm and its acceleration using our search algorithm (SRP-PHAT-DHBF). For the SRP-PHAT algorithm, we created a spatial energy map with the same resolution as the SRP-PHAT-DHBF algorithm at the finest level (10) of spatial subdivision, so that the map dimensions

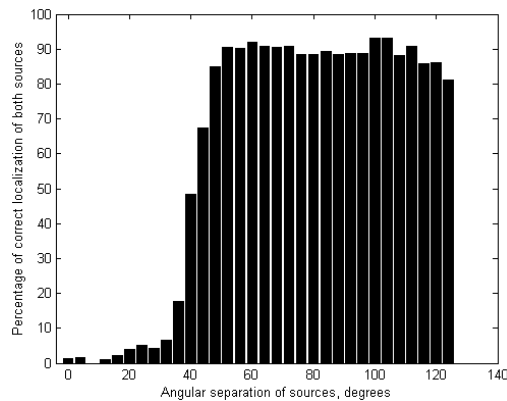


Fig. 5. Histogram of the percentage of correct localization of two simultaneous acoustic sources versus their angular separation for the SRP-PHAT-DHBF algorithm.

are 1024 by 1024 pixels and the resolution is about 0.18 degrees. A data frame length of 2048 points (93 ms at 22.05 kHz) was used. The difference in accuracy for the two algorithms is not very significant. In fact, SRP-PHAT-DHBF is even more accurate, perhaps by avoiding false peaks generated by SRP-PHAT at a fine resolution with noisy data. Other experiments conducted show that both SRP-PHAT-DHBF and SRP-PHAT exhibit the same level of robustness to noise in comparison with TDOA-based localizers.

In addition, the frequency hierarchy allows for selection of multiple sources at the coarsest level of the map as long as they are separated by certain minimal angular spacing that depends on array geometry and cutoff frequency. We performed tests to determine the minimum resolvable multiple source spacing when two sources were placed at random points (φ_1, θ_1) and (φ_2, θ_2) on the positive- Z hemisphere $\varphi \in [-\pi/2, \pi/2]$, $\theta \in [-\pi/2, \pi/2]$ of a fixed radius $\rho = 3.5m$ and the SRP-PHAT-DHBF algorithm was run to find out the probability of finding both sources as a function of angular distance between them. The histogram is presented in Figure 5 and shows that the spacing corresponding to 50% chance of correct separation is about 40 degrees.

D. Accuracy in real conditions

Finally, we performed a test of the algorithms in real reverberant and noisy conditions. We used a system consisting of the 0.3 m radius 7-channel microphone array, a data acquisition board and a PC for data collection. The experiments were conducted in a large office room $5.6 \times 2.7 \times 4.7$ (width \times height \times depth) meters. The microphone array configuration was the same as in the simulations above. We used Panasonic WM-61A omnidirectional speech-band button microphones with a custom preamplifier on AD797 chips connected to the PowerDAQ PD-MF-16-333/12L data acquisition board capturing 7 channels with 40.44 kHz sampling frequency and 12 bit resolution. The microphones were mounted on a large sheet of thick foam rubber to dampen wall reflections near the array, and the array was placed on one of walls, a third of

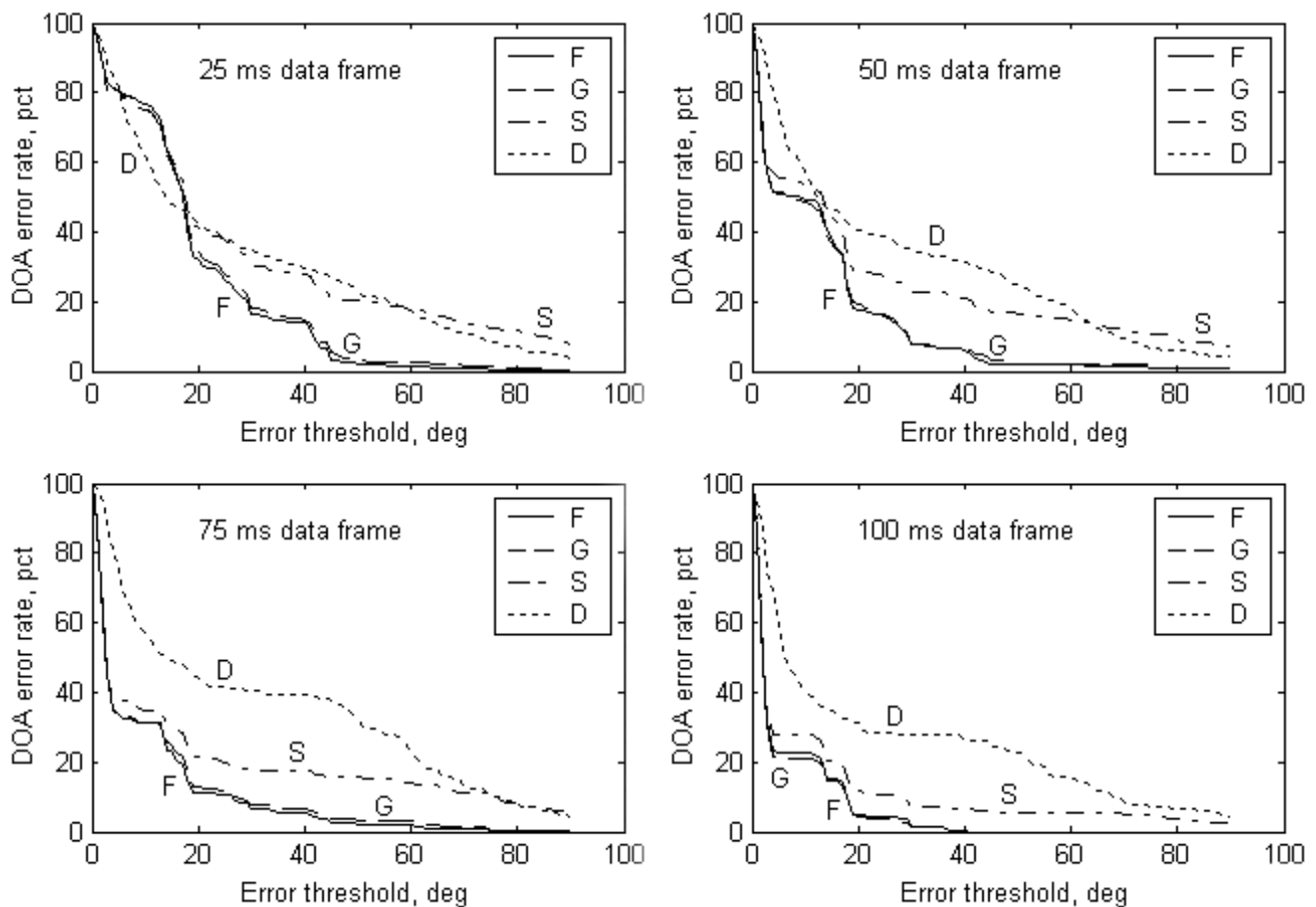


Fig. 6. Error rate (percents) versus error threshold (degrees) for varying frame length in real conditions.

the wall length from the corner. The measured reverberation time in the room was about 350 milliseconds, and the main source of noise in the room was the noise of computer equipment fans [25].

In the experiment, the person stood in five positions in the room and uttered the same sentence (“Dear fellow radio listeners!... Pass pass pass pass pass...”) at all positions. The positions were different in azimuth, elevation and distance, with the furthest one being about 2.5 meters from the array. From the recordings, we selected all frames where the SNR of the recorded speech was more than 12 dB and ran the same four algorithms on these frames. We repeated the processing with 4 different frame lengths (25, 50, 75 and 100 milliseconds). The total number of frames selected for the processing was 548, 276, 190 and 143 frames for frame length of 25, 50, 75 and 100 milliseconds, respectively. We show the results in Figure 6 using error rate versus error threshold metric (for example, if the dotted line shows an error rate of 70% at an error threshold of 7 degrees, it means that the estimate DOA produced by DHBF algorithm is within 7 degrees of true DOA in 30% of the cases, and the earlier lower values of error rate are achieved, the better the algorithm performance is).

It can be seen that the DHBF has a region of superior performance when the frame length is 25

milliseconds, and the performance is comparable to other algorithms using 50 milliseconds frames. When the frame length is 75 and 100 milliseconds, the performance of DHBF is notably worse than that of the stochastic region contraction and of the gradient descent. The negative effect of the short data frame on the full search is caused by the windowing operation which biases the bin phases and causes the false energy peak to appear (as described before, DHBF is often able to get a good estimation of DOA even when the full computed energy map has a maximum in an incorrect location because the bias is more pronounced at higher frequencies and the initial coarse-grid energy map estimation done by DHBF excludes them from consideration), and in real applications it is often necessary to track moving sound sources and to use short processing frames to minimize tracker latency, in which case DHBF is a viable alternative to existing fast localization algorithms.

E. Computational Speed

The main advantage of implementation of the proposed search algorithm is speed. In [4], the full spatial map with resolution of 0.1 degrees is computed in DOA-space for SRP and SRP-PHAT algorithms to compare their performance. In the proposed search algorithm, the number of evaluation of objective function is much smaller than in previously proposed fast search techniques. We performed a direct comparison between the number of operations needed to localize the source using full search, repeated gradient descent, stochastic region contraction and DHBF. We explicitly counted the number of evaluation of objective function in our code while doing source localization in real environment as described later on and found the following average number of evaluation of objective function per frame processed:

- Full search: $N_f = 2^{20} \approx 10^6$ evaluations;
- Repeated gradient descent: $N_g \approx 3.8 \cdot 10^4$ evaluations;
- Stochastic region contraction: $N_s \approx 4.6 \cdot 10^3$ evaluations;
- DHBF: $N_d \approx 3.7 \cdot 10^2$ evaluations. Further, most evaluations for DHBF are also *far cheaper* than the evaluations used in other algorithms (see below).

Indeed, for the maximum level of subdivision m , $k = 3$ sources and initial (coarse) grid subdivision level l , our algorithm performs $N_d = 2^{2l} + 4k(m - l)$ energy evaluations, compared to $N_f = 2^{2m}$ evaluations done if the full spatial map at fine resolution is computed. For $m = 10, l = 4, k = 3$ $N_d = 328$, which agrees with experimental results. Furthermore, most of the evaluations are performed with *bandpassed versions* of the signal, which additionally decreases the computational load for the DHBF. If we assume that the room is a cube with the side length of Z , the highest signal frequency is $f_d/2$ and frequency decimation is used according to the rules described above, then the number of evaluation weighted by the cost of each evaluation can be obtained by multiplying each evaluation by its weight $\xi(\lambda)$ which represents the ratio of the frequency range of the bandlimited signal to the frequency range of full-band

signal and is

$$N_{dd} = 2^{2l} \xi\left(\frac{Z}{2^l}\right) + 4k \sum_{k=l+1}^m \xi\left(\frac{Z}{2^k}\right), \quad (9)$$

where $\xi(\lambda)$ can be directly derived from the quadrant size heuristic as $\xi(\lambda) = c/(\lambda f_d)$ when $c < \lambda f_d$ and $\xi(\lambda) = 1$ otherwise. For the case considered above, using the same values of m, k, l , and $Z = 4m$ and $f_d/2 = 11$ kHz $N_{dd}/N_d \approx 0.19$, this results in an additional 5-fold reduction of the computational time.

We have implemented our algorithm together with several other several SRP-based algorithms and achieved real-time operation of DHBF on a dual Pentium III 600MHz Dell Precision 620 PC operating under Windows 2000 with no specialized hardware.

F. Vision-constrained source localization

One application area is source localization in multimodal user interface systems. In such system, *a priori* information from video will be available [26]. Knowing the location of the object in the image plane of the camera we can restrict the search to areas that are likely to be sources. In our implementation of this concept, an active camera is used to scan the room. A background model is constructed from several images taken at different camera orientations and mosaicing them. The room is constantly monitored for foreground objects, and a simple background subtraction method based on pixel intensities is used to classify every pixel as foreground/background. Since the camera image is 2D, the contour of the foreground object defines a “visual cone” in which the object lies, with the cone origin at the camera center. The cone is bounded by the room walls. The union of the visual cones, U , of all objects is either used directly in a full 3D search, by ignoring voxels in the initial coarse octree which do not intersect U or, for the 2D search in (ϕ, θ) space, the cone is reprojected back onto the (ϕ, θ) search space using known geometric relationships. In our implementation, we use two cameras, with one camera collecting the videoconferencing image from the active source, and the other camera constantly scanning the room, dynamically providing constraining data. In this way, system latency can be further decreased. Such uses are further described in [26].

V. CONCLUSIONS

We have presented a generic doubly hierarchical search algorithm for speech source localization using steered response power algorithms. The algorithm was designed with prior information about the speech in mind, and is able to achieve the accuracy that is comparable with other search-based steered response power algorithms in reduced processing time. Our experiments show that significant speedups can be achieved while keeping reasonable localization accuracy. The algorithm also has the ability to localize reasonably separated multiple active sources simultaneously. The search algorithm has applications in multiple areas including multimodal human-computer interaction, videoconferencing, and other entertainment, educational and remote collaboration applications.

REFERENCES

- [1] R. Duraiswami, D. N. Zotkin, L. S. Davis (2001). "Active speech source localization by a dual coarse-to-fine search", Proc. IEEE ICASSP 2001, Salt Lake City, UT, vol. 5, pp. 3309-3312.
- [2] M. S. Brandstein (1999). "Time-delay estimation of reverberated speech exploiting harmonic structure", J. Acoust. Soc. Am., vol. 105, no. 5, pp. 2914-2919.
- [3] M. S. Brandstein and D. B. Ward (editors) (2001). "Microphone Arrays: Signal Processing Techniques and Applications", Springer-Verlag, Berlin, Germany.
- [4] J. H. DiBiase, H. F. Silverman, and M. S. Brandstein (2001). "Robust localization in reverberant rooms", in Microphone Arrays: Signal Processing Techniques and Applications, ed. by M. S. Brandstein and D. B. Ward, Springer-Verlag, Berlin, Germany, pp. 157-180.
- [5] D. Colton and R. Kress (1992). "Inverse acoustic and electromagnetic scattering theory", Springer-Verlag, Berlin, Germany.
- [6] M. S. Brandstein and S. Griebel (2001). "Explicit speech modeling for microphone array applications", in Microphone Arrays: Signal Processing Techniques and Applications, ed. by M. S. Brandstein and D. B. Ward, Springer-Verlag, Berlin, Germany, pp. 131-153.
- [7] M. Wax and T. Kailath (1983). "Optimum localization of multiple sources by passive arrays", IEEE Transactions on Acoustics, Speech and Signal Processing, vol. 31, pp. 1210-1218.
- [8] V. M. Alvarado (1990). "Talker localization and optimal placement of microphones for a linear microphone array using stochastic region contraction", Ph.D. Thesis, Brown University, Providence, RI, USA.
- [9] M. F. Berger and H. F. Silverman (1991). "Microphone array optimization by stochastic region contraction", IEEE Transactions on Signal Processing, vol. 39, no. 11, pp. 2377-2386.
- [10] A. Doucet, N. de Freitas, and N. Gordon (editors) (2001). "Sequential Monte Carlo methods in practice", Springer-Verlag, Berlin.
- [11] J. Vermaak and A. Blake (2001). "Nonlinear filtering for speaker tracking in noisy and reverberant environments", Proc. IEEE ICASSP 2001, Salt Lake City, UT, vol. 5, pp. 3021-3024.
- [12] D. B. Ward and R. C. Williamson (2002). "Particle filter beamforming for acoustic source localization in a reverberant environment", Proc. IEEE ICASSP 2002, Orlando, FL, vol. 2, pp. 1777-1780.
- [13] H. Silverman (1987). "Some analysis of microphone arrays for speech data acquisition", IEEE Transactions on Acoustics, Speech and Signal Processing, vol. 35, no. 12, pp. 1699-1712.
- [14] W. M. Hartmann (1999). "How We Localize Sound," Physics Today, vol. 52, no. 11, pp. 24-29.
- [15] C. H. Knapp and G. C. Carter (1976). "The generalized correlation method for estimation of time delay", IEEE Transactions on Acoustics, Speech and Signal Processing, vol. 24, no. 4, pp. 320-327.
- [16] M. S. Brandstein and H. F. Silverman (1997). "A robust method for speech signal time-delay estimation in reverberant rooms", Proc. IEEE ICASSP 1997, Munich, Germany, pp. 375-378.
- [17] M. Omologo and P. Svaizer (1997). "Use of the crosspower-spectrum phase in acoustic event location", IEEE Transactions on Speech and Audio Processing, vol. 5, no. 12, pp. 288-292.
- [18] H. Wang and P. Chu (1997). "Voice source localization for automatic camera pointing system in videoconferencing", Proc. IEEE ICASSP 1997, Munich, Germany, pp. 187-190.
- [19] J. H. DiBiase (2000). "A high-accuracy, low-latency technique for talker localization in reverberant environments", Ph.D. Thesis, Brown University, Providence, RI, USA.
- [20] P. P. Vaidyanathan (1992). "Multirate systems and filter banks", Prentice Hall, Upper Saddle River, NJ.
- [21] H. Samet (1990). "Applications of spatial data structures", Addison-Wesley, Boston, MA.
- [22] M. S. Brandstein and H. F. Silverman (1997). "A practical methodology for speech source localization with microphone arrays", Computer, Speech and Language, vol. 11, no. 2, pp. 91-126.
- [23] J. B. Allen and D. A. Berkeley (1979). "Image method for efficiently simulating small-room acoustics", J. Acoust. Soc. Am., vol. 65, no. 5, pp. 943-950.
- [24] R. Duraiswami, N. A. Gumerov, D. N. Zotkin, L. S. Davis (2001). "Efficient evaluation of reverberant sound fields", Proc. IEEE WASPAA 2001, New Paltz, NY, October 2001, pp. 203-206.
- [25] D. N. Zotkin, R. Duraiswami, L. S. Davis, and I. Haritaoglu (2000). "An audio-video front end for multimedia applications", Proc. IEEE SMC 2000, Nashville, TN, pp. 786-791.
- [26] D. N. Zotkin, R. Duraiswami, V. Philomin, and L. S. Davis (2000). "Smart videoconferencing", Proc. IEEE ICME 2000, New York City, NY, vol. 3, pp. 1597-1600.

VI. LIST OF FIGURES

Figure 1. Beamformer peak width as a function of frequency.

Figure 2. a) Sample energy map (energy vs azimuth and elevation) obtained using SRP-PHAT algorithm with three sources, resolution 512x512. b) Same map for the signals lowpassed at 2 kHz. c) Same map for the signals lowpassed at 2 kHz, resolution 16x16.

Figure 3. a) Average localization error for 4 algorithms. b) Percentage of correct localization for 4 algorithms. c) Average localization error with data frames where full search fails omitted. d) Percentage of correct localization under the same condition. R is the reverberation time in milliseconds.

Figure 4. Localization error versus SNR for SRP-PHAT and SRP-PHAT-DHBF.

Figure 5. Histogram of the percentage of correct localization of two simultaneous acoustic sources versus their angular separation for the SRP-PHAT-DHBF algorithm.

Figure 6. Error rate (percents) versus error threshold (degrees) for varying frame length in real conditions.