

# MUSIC 法を用いた音源定位のベイズ拡張

大塚 琢馬<sup>1</sup>, 中臺 一博<sup>2,3</sup>, 尾形 哲也<sup>1</sup>, 奥乃 博<sup>1</sup>

<sup>1</sup> 京都大学大学院 情報学研究科 知能情報学専攻

<sup>2</sup> (株) ホンダ・リサーチ・インスティテュート・ジャパン

<sup>3</sup> 東京工業大学大学院 情報理工学研究科 情報環境学専攻

## 1. はじめに

音響情報は人間の知覚の重要な位置を占める。例えば、人は足音を聞くことで目に頼ることなく誰かが近づいてきている、あるいは遠ざかっているといった状況を理解することができる。ロボットや計算機による周囲の音響情報の理解、つまり、「音環境理解」の実現は、聴覚障害者の補助や、人間の音に対する気づきを向上させることができると期待される [1]。

音源定位はマイクロフォンアレイを用いた同時発話混合音声の分離 [2]、遠隔ロボットのオペレータへの音源方向提示 [3]、移動ロボットによる音源検出と位置推定 [4] など、音環境理解にとって重要な要素技術である。図 1 に示すような、複数音源、ロボットの移動、音源移動など、動的に音環境が変化する状況においても、手間のかかるパラメータ設定をしなくてもロボットが頑健に各音源を定位、追跡することが望まれる。

マイクロフォンアレイを用いた音源定位法はビームフォーミングに基づく手法 [5] と、MUltiple SIgnal Classification (MUSIC) に基づく手法 [6, 7, 8] がロボットによく応用される。我々は次の理由より、MUSIC 法を利用する。(1) MUSIC の方が雑音に頑健である、(2) 音源数がマイクロフォン数未満という条件下では、比較的安定して複数音源の定位が可能である。

通常の MUSIC 法では、音源が到来しているかどうかを MUSIC スペクトルと呼ばれる音源到来評価関数に対して閾値を設定して判定する。多くの場合、適切な閾値は環境中の音源数や残響時間などに依存するため、手動設定が困難である。MUSIC 法を用いた場合の環境中の音源数推定問題は、赤池情報量規準の利用 [8] や、サポートベクターマシンの適用 [9] によってこれまで取り組まれてきた。しかし、これらの手法で音源数が推定できたとしても、適切な音源検出閾値を設定するという問題は依然として残っている。この問題に対する典型的な対策としては、マイクロフォンアレイを設置した環境で録音した音響信号から計算した MUSIC スペクトルを見ながら手動で閾値を設定するという方法であった。

本稿では、MUSIC 法による音源定位のベイズ拡張を行い、従来法で必要とされた閾値に相当する情報を自動的に学習することを試みる。これにより、閾値設定の手間を省くと共に、試行錯誤により設定された閾値の精度と同等以上の定位精度を実現する。本手法は次の 2 つのステップから成る。(1) マイクロフォンアレイが置かれた環境で録音した数十秒程度の音響信号から、音源存在閾値に相当するパラメータを学習する。学習には変分ベイズ隠れマルコフモデル (VB-HMM) [10] に基づくパラメータ推定アルゴリズムを用いる。(2) VB-HMM により学習したパラメータを用いた複数音源の逐次的定位を行う。逐次定位では、観測モデルが VB-HMM より複雑になるため、パーティクルフィルタ [11] を用いる。

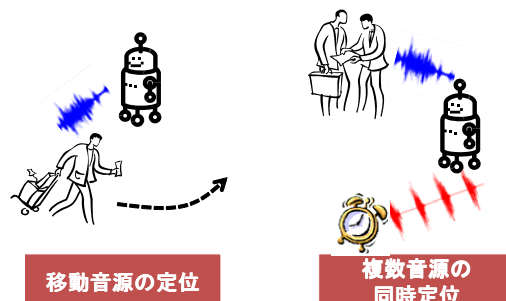


図 1 動的環境下での音源定位

## 2. MUSIC 法を用いる音源定位

まず本稿が扱う問題を述べ、MUSIC スペクトルの算出法を説明する。本稿での水平面上の音源到来方向推定問題を、図 2 に示した。今回用いたマイクロフォンアレイは、マイクロフォンがロボットに円状に 8 本配置されており、水平面上に  $5^\circ$  刻みの解像度での定位を行う。以下に本稿で扱う問題設定を示す。

入力  $M$  チャンネルの音響信号と、各周波数ビンごとに  $D$  方向からの伝達関数、  
出力  $N$  個の音源到来方向、  
仮定 同時に検出可能な最大音源数  $N_{max}$  はマイクロフォンの数未満 ( $N \leq N_{max} < M$ )。

水平面一周を  $5^\circ$  刻みに定位するので、 $D = 72$  である。

次に、MUSIC スペクトルの算出法について簡単に述べる。より詳細は文献 [6, 8] などに記述されている。MUSIC 法は時間周波数領域<sup>1</sup>において適用される。

$\mathbf{x}_{\tau,\omega} \in \mathbb{C}^M$  を  $M$  チャンネル音響信号の時間フレーム  $\tau$ 、周波数ビン  $\omega$  における複素振幅ベクトルとする。各周波数ビン  $\omega$ 、 $\Delta T$  [sec] 間隔の時刻  $t$  に対して、(1) 入力信号の自己相関行列  $\mathbf{R}_{t,\omega}$  の計算、(2)  $\mathbf{R}_{t,\omega}$  の固有値分解、(3) 固有ベクトルと伝達関数を用いた MUSIC スペクトルの計算を行う。

(1) 入力信号の自己相関行列は時間  $\Delta T$  で観測したサンプル値の相関として計算する。

$$\mathbf{R}_{t,\omega} = \frac{1}{\hat{\tau}(t) - \hat{\tau}(t - \Delta T)} \sum_{\tau=\hat{\tau}(t-\Delta T)}^{\hat{\tau}(t)} \mathbf{x}_{\tau,\omega} \mathbf{x}_{\tau,\omega}^H, \quad (1)$$

ただし、 $(\cdot)^H$  はエルミート転置、 $\hat{\tau}(t)$  は時刻  $t$  に対応する時間フレームを表す。入力ベクトル  $\mathbf{x}_{\tau,\omega}$  の  $M$  個の要素は各チャンネルに対応する。

(2)  $\mathbf{R}_{t,\omega}$  を次のように固有値分解する。

$$\mathbf{R}_{t,\omega} = \mathbf{E}_{t,\omega} \mathbf{Q}_{t,\omega} \mathbf{E}_{t,\omega}^H, \quad (2)$$

ここで、 $\mathbf{E}_{t,\omega}$  は固有ベクトル、 $\mathbf{Q}_{t,\omega}$  は固有値から成る対角行列である。 $\mathbf{E}_{t,\omega} = [\mathbf{e}_{t,\omega}^1 \dots \mathbf{e}_{t,\omega}^M]$  と、 $\mathbf{R}_{t,\omega}$  の  $M$  個の固

<sup>1</sup>我々の実装では、サンプリング周波数 16000 [Hz] で、窓長 512 [pt]、シフト幅 160 [pt] の短時間フーリエ変換を行っている。



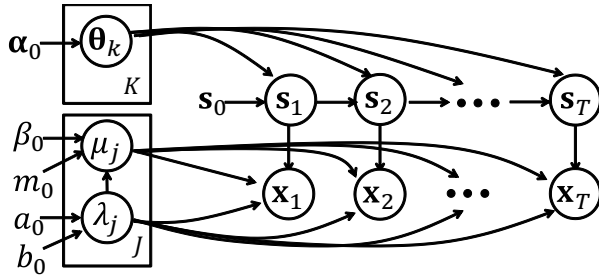


図4 VB-HMMのグラフィカルモデル

表1 隣接状態を考慮した状態遷移の場合分け

前状態 $s_{t-1,d}$	隣接前状態 $1 - s_{t-1,d-1} - s_{t-1,d+1}$	音源存在確率 $p(s_{t,d} = 1   s_{t-1,d-1:d+1})$
0 (off)	0	$\theta_1$
0 (off)	1	$\theta_2$
1 (on)	0	$\theta_3$
1 (on)	1	$\theta_4$

から、次状態で音源が出現する、継続する、消滅するといった遷移を考える。本稿ではさらに、移動する音源についても考慮するために、表1のように前状態の組み合わせから成る4つの場合を考える。すなわち、前時刻の同方向ビン  $s_{t-1,d}$  に音源が存在するかどうか、前時刻の隣接方向ビン  $s_{t-1,d\pm 1}$  のいずれかに音源が存在するかによって分類する。例えば、 $\theta_1$  は前時刻に当該方向  $d$  及び隣接ビン  $d \pm 1$  に音源が存在しない状態から音源が出現する確率、 $\theta_2$  は、前時刻に方向  $d$  に音源が存在しないが、隣接ビン  $d \pm 1$  には音源が存在したため、その音源が方向  $d$  に移動してきて  $s_{t,d} = 1$  となる確率を表す。状態遷移確率は以下の通り。

$$p(s_t | s_{t-1}, \theta) = \prod_{d=1}^D \prod_{k=1}^4 \prod_{j=0}^1 (\theta_k^{s_{t,d}} (1 - \theta_k)^{1-s_{t,d}})^{f_k(s_{t-1,d})} \quad (8)$$

ここで、 $f_k(s_{t-1,d})$  は表1に従って、方向ビン  $d$  の周りの前状態の値  $s_{t-1,d-1}, s_{t-1,d}, s_{t-1,d+1}$  によって条件  $k$  に合致するときに  $f_k(\cdot, d) = 1$  その他の場合は0を返す条件識別関数である。初期状態としては、音源は存在しない、すなわちすべての  $d$  に対して  $s_{0,d} = 0$  とする。

状態遷移パラメータである  $\theta = [\theta_1, \dots, \theta_4]$  には、式(8)の共役事前分布としてベータ分布を用いる。

$$p(\theta | \alpha_0) = \prod_{k=1}^4 \mathcal{B}(\theta_k | \alpha_{0,1}, \alpha_{0,2}), \quad (9)$$

ただし、 $\mathcal{B}(\cdot | c, d)$  はパラメータ  $c, d$  を持つベータ分布の確率密度関数である。

### 3.1.3 事後分布の推定

VB-HMMの学習は、事後分布  $p(s_{1:T}, \theta, \mu, \lambda | \mathbf{x}_{1:T})$  を以下のように因数分解可能な分布に近似して推定する。

$$p(s_{1:T}, \theta, \mu, \lambda | \mathbf{x}_{1:T}) \approx q(s_{1:T}, \theta, \mu, \lambda), \quad (10)$$

$(\cdot)_{1:T}$  は、時刻1から  $T$  までの確率変数の集合を表す。文献[10]に一般的なVB-HMMの推論について述べられている。以下では、分布がどのように更新されるかを簡単に述べる。 $q(\theta) = \prod_k q(\theta_k)$  はそれぞれの  $k$  に対し、式(11)に示すパラメータ  $\hat{\alpha}_{k,0}, \hat{\alpha}_{k,1}$  を持つベータ分布となり、 $q(\mu, \lambda) = \prod_j q(\mu_j, \lambda_j)$  は、式(12), (13)のように、パラメータ  $\hat{\beta}_j, \hat{m}_j, \hat{a}_j, \hat{b}_j$  を持つ正規ガウス分布となる。

$$\hat{\alpha}_{k,j} = \alpha_{0,j} + \sum_{t,d} \langle s_{t,d,j} f_k(s_{t-1,d}) \rangle, \quad (11)$$

$$\hat{\beta}_j = \beta_0 + w_j, \hat{m}_j = (\beta_0 m_0 + w_j \bar{x}_j) / (\beta_0 + w_j), \quad (12)$$

$$\hat{a}_j = a_0 + \frac{w_j}{2}, \hat{b}_j = b_0 + \frac{w_j S_j^2}{2} + \frac{\beta_0 w_j (\bar{x}_j - m_0)^2}{2(\beta_0 + w_j)}, \quad (13)$$

ただし、変数  $s_{t,d,j}$  は、 $s_{t,d} = 0$  のとき、 $s_{t,d,0} = 1$ 、また、 $s_{t,d} = 1$  のとき、 $s_{t,d,1} = 1$  となる変数である。式(12), (13)に用いられる正規分布の十分統計量は

$$w_j = \sum_{t,d} \langle s_{t,d,j} \rangle, \bar{x}_j = \frac{\sum_{t,d} \langle s_{t,d,j} x_{t,d} \rangle}{w_j}, S_j^2 = \frac{\sum_{t,d} \langle s_{t,d,j} (x_{t,d} - \bar{x}_j)^2 \rangle}{w_j}.$$

と定義する。また、 $\langle \cdot \rangle$  は式(10)の分布による期待値演算子である。各時刻の状態変数と状態遷移の期待値  $\langle s_{t,d,j} \rangle$ ,  $\langle s_{t,d,j} f_k(s_{t-1,d}) \rangle$  は次のように計算する。

$$\langle s_{t,d,j} \rangle \propto \alpha(s_{t,d,j}) \beta(s_{t,d,j}), \quad (14)$$

$$\langle s_{t,d,j} f_k(s_{t-1,d}) \rangle \propto \tilde{\alpha}(s_{t-1,d,k}) \tilde{p}(s_{t,d} | s_{t-1}) \tilde{p}(x_{t,d} | s_{t,d}) \beta(s_{t,d,j}), \quad (15)$$

ただし、 $\alpha(s_{t,d,j})$  と  $\beta(s_{t,d,j})$  はそれぞれ前向き・後ろ向き再帰式により計算される。

$$\alpha(s_{t,d,j}) \propto \sum_{k=1}^4 \tilde{\alpha}(s_{t-1,d,k}) \tilde{p}(s_{t,d} | s_{t-1}) \tilde{p}(x_{t,d} | s_{t,d}), \quad (16)$$

$$\beta(s_{t,d,j}) = \sum_{j'=0}^1 \beta(s_{t+1,d,j'}) \tilde{p}(s_{t+1,d,j'} | s_{t,d,j}) \tilde{p}(x_{t,d} | s_{t,d}), \quad (17)$$

式(14), (15) 遷移、観測確率の幾何平均は次の通り。

$$\tilde{p}(s_{t,d} = j | s_{t-1}) \propto \exp \{ \psi(\hat{\alpha}_{k,j}) - \psi(\hat{\alpha}_{k,0} + \hat{\alpha}_{k,1}) \}, \quad (18)$$

$$\tilde{p}(x_{t,d} | s_{t,d}) \propto \prod_j \exp \left\{ \frac{\psi(\hat{a}_j) - \log \hat{b}_j - 1/\hat{\beta}_j}{2} - \frac{a_j (x_{t,d} - \hat{m}_j)^2}{2\hat{b}_j} \right\}^{s_{t,d,j}} \quad (19)$$

式(14), (15) はともに、添字  $j, k$  を動かしたとき総和が1になるように正規化されている。 $\tilde{\alpha}(s_{t-1,d,k})$  は、状態遷移の条件  $k$  に関する前向き確率である。本節で示されたパラメータ更新式(11)–(15)が収束するまで計算される。初期値としては、 $\langle s_{t,d,j} \rangle$  と  $\langle s_{t,d,j} f_k(s_{t-1,d}) \rangle$  の値を、観測変数  $x_{t,d}$  の値を  $m_0$  の値を閾値として処理することで、0ないし1を与える。

### 3.2 パーティクルフィルタによるオンライン音源定位

本節ではパーティクルフィルタ[11]を用いた、オンライン音源定位手法を述べる。オンライン推定では、式(11)–(13)で求めたパラメータの事後分布を利用する。パーティクルフィルタの推定対象は、MUSICスペクトルの時系列データが与えられたときの、各方向ビンにおける音源存在事後確率である。この分布を  $P$  個のパーティクルを用いて以下のように近似計算する。

$$p(s_t | \mathbf{x}_{1:T}) \approx w_p s_t^p, \quad (20)$$

ただし、 $w_p$  はパーティクル  $p$  の重み、 $s_t^p$  は状態ベクトルの値である。これらの  $w_p$  と  $s_t^p$  は次のように得る。

(1) 提案分布から  $s_t^p$  をサンプルする。

$$s_t^p \sim q(s_t | \mathbf{x}_t, m, a, b), \quad (21)$$

$$q(s_t^p | \mathbf{x}_t, \hat{m}, \hat{a}, \hat{b}) \propto \prod_d C(x_{t,d}) \exp(-\Delta_{d,j}^2 / 2) \hat{a}_j^{s_{t,d,j}}, \quad (22)$$

ただし、 $x_{t,d}$  が極大値を取る  $d$  のとき、 $C(x_{t,d}) = 1$  でその他の場合は  $C(x_{t,d}) = 0$  となる。提案分布の重みにはマハラノビス距離  $\Delta_{d,j}^2 = (x_{t,d} - \hat{m}_j)^2 \hat{a}_j / \hat{b}_j$  を用いる。

(2) 各パーティクル  $p$  について、重み  $w_p$  を算出。

$$w_p \propto \frac{\tilde{p}(\mathbf{x}_t | s_t^p) \tilde{p}(s_t^p | s_{t-1}^p)}{q(s_t^p | \mathbf{x}_t, \hat{m}, \hat{a}, \hat{b})}, \quad (23)$$

$$\tilde{p}(\mathbf{x}_t | s_t^p) = \int p(\mathbf{x}_t | s_t^p, \mu, \lambda) q(\mu, \lambda) d\mu d\lambda, \quad (24)$$

$$\tilde{p}(s_t^p | s_{t-1}^p) = \int p(s_t^p | s_{t-1}^p, \theta) q(\theta). \quad (25)$$

<sup>2</sup> $\psi(\cdot)$  はディガンマ関数。

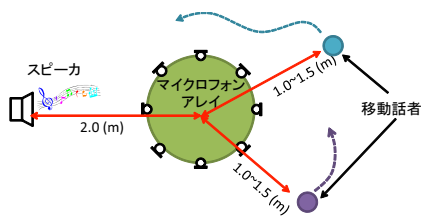


図5 実験条件: マイクアレイの周囲を動く移動話者と固定音源

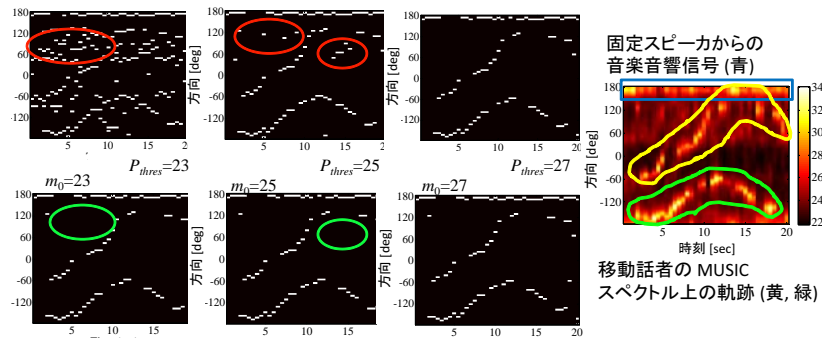


図6 音源定位結果: 白が音源が存在する方向, 時間ピン. 上図: 固定閾値  $P_{thres}$  による定位結果. 下図: 初期値  $m_0$  を用いた本手法の定位. 右図: 観測された対数 MUSIC スペクトル. 音楽音響信号が 180 [deg] 付近に存在し, 2 人の話者が移動している.

式 (24),(25) にある状態遷移, 観測確率は, VB-HMM で計算された式 (6),(8) の事後分布で積分消去することで計算できる. 分布の共役性を用いると, この積分計算は次のように解析的に求まる.

$$\bar{p}(\mathbf{x}_t | \mathbf{s}_t^p) = \prod_d St(x_{t,d} | \hat{m}_j, \frac{\hat{\beta}_j \hat{a}_j}{(1 + \hat{\beta}_j) \hat{b}_j}, 2\hat{a}_j) s_{t,d}^p, \quad (26)$$

$$\bar{p}(\mathbf{s}_t^p | \mathbf{s}_{t-1}^p) = \prod_d \prod_k \left( \frac{\hat{\alpha}_{k,s_{t,d}}}{(\hat{\alpha}_{k,0} + \hat{\alpha}_{k,1})} \right)^{f_k(s_{t-1}^p, d)} \quad (27)$$

ただし,  $St(\cdot | m, \lambda, \nu)$  は平均  $m$ , 精度  $\lambda$ , 自由度  $\nu$  の Student t-分布である. さらに, 最大の音源数を  $N_{max}$  に抑えるため, 状態ベクトル  $\mathbf{s}_t^p$  に存在する音源数が  $N_{max}$  を超える場合には観測確率は 0 とする.

全パーティクルの重み計算後, 各パーティクルの重み  $w_p$  は  $\sum_{p=1}^P w_p = 1$  となるよう正規化する. この手順に従い, 式 (20) の音源存在の事後分布を計算する. 我々の実装手法では, 各ステップごとにパーティクルが持つ重みに比例してリサンプリング処理が行われる.

#### 4. 評価実験

評価実験では, VB-HMM によるパラメータ分布推定とパーティクルフィルタを用いたオンライン音源定位から成る本手法と, 従来の固定閾値を用いて音源定位する手法を比較する. オフラインでの VB-HMM での学習は, 1 人の話者がマイクアレイの周囲を発話しながら動く音響信号で行った. オンラインの音源定位実験に使用した音源の配置を図 5 に示す. マイクアレイの周囲を移動する 2 話者と, 固定されたスピーカから音楽が再生されている. オフライン, オンラインで用いられた信号の長さはともに 20 [sec] である. パラメータの設定は次の通り.  $N_{max} = 3$ ,  $\alpha_0 = [1, 1]$ ,  $\beta_0 = 1$ ,  $a_0 = 1$ ,  $b_0 = 500$ . パーティクル数は  $P = 500$  とした. 実験で使った室内の残響時間は  $RT_{20} = 840$  [msec] であった.

図 6 にオンライン音源定位の結果を示す. 従来法の閾値は  $P_{thres} = 23, 25, 27$  に設定されており, 本手法の初期値は  $m_0 = 23, 25, 27$  に設定されている. パーティクルフィルタの定位結果の図では, 事後分布の音源存在確率が 0.95 以上のピンを音源が存在するとして白く表示している. 従来法においては, 閾値を低く設定した場合は図 6 の赤枠で示すように音源の誤検出が頻発する. 対して, 本手法では緑枠で示すように, 学習の初期値に対して頑健に妥当な音源定位結果を示している. また, 本手法において音源存在確率の閾値を 0.95-1.00 まで動かして結果を検証したが, これらの値を閾値に対しても頑健に同様の結果を示すことを確認した. この結果から, 本手法におけるオフライン学習, オンライン定位の枠組

みが, 自動的に音源定位に適したパラメータに収束することが確認できる. さらに, 今回の実験条件から, 本手法は学習時に 1 音源しか用いなくても, 複数音源も安定してオンライン定位することが実証された.

#### 5. まとめと今後の課題

本稿では MUSIC 法に基づく音源定位法のベイズ拡張を述べた. 本手法は, (1) VB-HMM によるパラメータの自動学習, (2) パーティクルフィルタを用いたオンライン音源定位から成る. 評価実験では,  $RT_{20} = 840$  [msec] の残響環境下で, 1 音源の音響信号の学習に対し, 3 音源同時音源定位を実現した. 今後の課題としては, 音源定位の時系列や, 各音源の音色などを考慮した音源トラッキングへの拡張や, ロボット聴覚システム HARK [2] への実装, また, 実際に移動ロボットに本手法を適用して, ロボット位置と環境中に存在する音源位置の推定を通じた音環境理解システムの構築などが挙げられる.

謝辞: 本研究の一部は科研費特別研究員奨励金/基盤 (S), JST-ANR BINAHR, GCOE の支援を受けた.

#### 参考文献

- [1] Y. Kubota et al., "Design and Implementation of 3D Auditory Scene Visualizer towards Auditory Awareness with Face Tracking," in *Proc. of IEEE Int'l Symposium on Multimedia (ISM-2008)*, 2008, pp. 468-476.
- [2] K. Nakadai et al., "Design and Implementation of Robot Audition System "HARK"," *Advanced Robotics*, vol. 24, no. 5-6, pp. 739-761, 2010.
- [3] T. Mizumoto et al., "Design and Implementation of Selectable Sound Separation on a Texai Telepresence System using HARK," in *Proc. of ICRA*, 2011, pp. 2130-2137.
- [4] Y. Sasaki et al., "Map-Generation and Identification of Multiple Sound Sources from Robot in Motion," in *Proc. of IROS*, 2010, pp. 437-443.
- [5] S. Doclo and M. Moonen, *Microphone arrays*. Springer, 2001, ch. GSVD-based optimal filtering for multi-microphone speech enhancement, pp. 111-132.
- [6] R. O. Schmidt, "Multiple Emitter Location and Signal Parameter Estimation," *IEEE Trans. on Antennas and Propagation*, vol. 34, no. 3, pp. 276-280, 1986.
- [7] F. Asano et al., "Real-time Sound Source Localization and Separation System and Its Application to Automatic Speech Recognition," in *Proc. of Eurospeech*, 2001, pp. 1013-1016.
- [8] P. Danès and J. Bonnal, "Information-Theoretic Detection of Broadband Sources in a Coherent Beam-space MUSIC Scheme," in *Proc. of IROS*, 2010, pp. 1976-1981.
- [9] K. Yamamoto et al., "Detection of Overlapping Speech in Meeting using Support Vector Machines and Support Vector Regression," *IEICE Trans. Fundamentals*, vol. E89-A, no. 8, pp. 2158-2165, 2006.
- [10] M. J. Beal, "Variational Algorithms for Approximate Bayesian Inference," Ph.D. dissertation, Gatsby Computational Neuroscience U., Univ. Colledge London, 2003.
- [11] M. Arulampalam et al., "A Tutorial on Particle Filters for Online Nonlinear/Non-Gaussian Bayesian Tracking," *IEEE Trans. on Signal Proc.*, vol. 50, no. 2, pp. 174-189, 2002.