

楽曲中の歌声の基本周波数と音素を同時推定可能なフレームワーク*

©藤原弘将（産総研、京大）、後藤真孝（産総研）、奥乃博（京大）

1 はじめに

音楽は、産業的にも文化的にも重要なコンテンツであり、その中でも歌声は重要な役割を果たしている。本稿では、混合音中の歌声の歌詞（音素）と基本周波数（F0）を同時に認識するための手法、W-PST（Weighted composition of Probabilistic Spectral Template）法を提案し、F0推定と音素認識の実験によりその有効性を確認する。本稿では歌詞とF0についてのみ触れるが、提案する手法は声質（歌手名）など歌声のその他の要素の認識にも適用可能であり、混合音中の歌声を扱うための新たなフレームワークと位置づけることができる。

歌詞は歌手が歌声によって伝えたい内容を表現し、F0は楽曲の旋律を表すと同時に、歌手の技巧や表情なども表現するため、どちらも歌声を構成する重要な要素である。そのため、混合音中からこれらの要素を自動認識する技術は、音楽情報検索などにも応用可能で、重要な基礎技術となる。しかし、歌声は話し声に比べて、ビブラートやF0の変化幅の広さ、歌手の感情表現などに起因する変動が多い上に、伴奏音が大音量で重畳するため、歌声（音素）の自動認識は非常に難しい研究課題である。

我々は、今までに音楽と歌詞の時間的対応付け手法[1, 2]と混合音中の歌声のF0推定手法[3]について研究してきた。これらの手法では共通して、混合音から調波構造を手がかりに音を分離し、それを統計的手法により識別するというアプローチをとっていた。具体的には、歌詞の時間的対応付けの場合、既存手法によって推定された歌声のF0の音がどの音素であるかを識別し、歌声のF0推定の場合、各時刻の周波数成分の候補が歌声であるかそれ以外の音であるかを識別していた。

しかし、それらの手法は下記の2つの問題点を抱えていた。

分離の問題 歌声の認識性能が、その前段に行われる分離の性能に大きく依存していた。そのため、F0推定や、分離の際にスペクトルから調波成分を選択する処理のエラーが、性能に悪影響を与えていた。

スペクトル包絡推定の問題 従来の我々の手法では、スペクトル包絡を分離後の歌声の調波構造から推定しスペクトル包絡同士の距離を計算することで、歌声を認識していた。しかし、与えられた調波構造から元のスペクトル包絡を一意に復元することは原理的に不可能であり、F0が高い音など、調波構造の各倍音成分の谷間の幅が広い場合など、距離を正確に計算することが困難であった。

混合音中の歌詞または音素の認識に関する関連研

究として、[4, 5, 6, 7]がある。いずれの研究も、歌声を分離しているか、もしくは、そもそも伴奏の影響を考慮していないかで、この述べた問題は解決されていなかった。混合音中の歌声に対するF0推定の研究として、[8, 9]がある¹が、本研究のように歌声のスペクトル包絡をモデル化し学習することで歌声のF0を推定しているものはなかった。

本稿では、これらの問題点を解決する新しい手法を提案する。この手法は、歌声を分離したり、単一の調波構造からスペクトル包絡を推定したりせず、観測されたスペクトルを伴奏音が重畳したありのままの形を確率的にモデリングする。さらに、学習の過程では、複数の調波構造を用いることで、より正確にスペクトル包絡を推定する。

2 歌声を認識するための新たなフレームワーク

図1(c)と(d)で示されるように、歌声を含む混合音のスペクトルがある確率分布の集合から生成されると仮定する。本稿では、それを**確率的スペクトルテンプレート**（Probabilistic Spectral Template）と呼ぶ。ここで、スペクトルの各ビンのパワーはある確率分布に従い、その確率分布はスペクトルのビンごとに異なると考える。スペクトルの加法性を仮定すると、確率的スペクトルテンプレートは、歌声を表現するスペクトルテンプレート（図1(a)）と歌声以外の音を表現するスペクトルテンプレート（図1(b)）の線形軸上での加算で表現することができる。前者を**歌声スペクトルテンプレート**（Vocal Spectral Template）、後者を**ノイズスペクトルテンプレート**（Noise Spectral Template）と呼ぶ。それらの2つのスペクトルテンプレートの加算の際に重みパラメータを導入し、重み付きで加算することで、様々なS/N比のスペクトルを表現できる。さらに、ソースフィルターモデルを仮定すると、歌声スペクトルテンプレートは、スペクトル包絡を表現する**歌声包絡テンプレート**（Vocal Envelope Template）（図2(a)）と駆動源の調波構造を表現する**調波フィルタ**（Harmonic Filter）（図2(b)）の積によって生成されると考えられる。調波フィルタの形状は、F0の値をパラメータとして、コントロールできる。ここで、この確率モデルのパラメータである調波フィルタのF0と、歌声・ノイズスペクトルテンプレートのそれぞれの重みが定まれば、観測スペクトルのモデルに対する尤度を計算することができる。このモデルを用いると、各音素を表現する歌声包絡テンプレートをあらかじめ学習しておき観測スペク

¹歌声に限定しない一般のメロディのF0推定の研究は数多くあるが、ここでは歌声に特化したもののみを紹介する。

*A framework for concurrently estimating F0 and phonemes of singing voice in music. by FUJIHARA, Hiromasa (AIST, Kyoto University), GOTO, Masataka (AIST), and OKUNO, G. Hiroshi (Kyoto University)

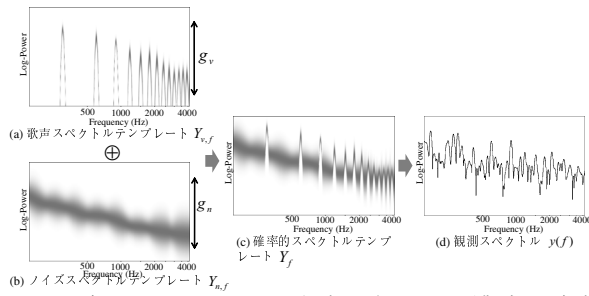


Fig. 1. 観測スペクトルの生成過程. 図の濃淡は確率密度を表現する. 重みパラメータ g_v と g_n を調整することで, 様々な S/N 比のスペクトルを表現できる.

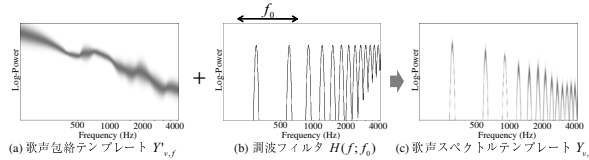


Fig. 2. 歌声スペクトルテンプレートの例. 歌声包絡テンプレートと調波フィルタから生成される.

トルに対して最尤な歌声包絡テンプレートを選択することで音素認識ができ, 最尤な F0 の値を推定することで F0 推定ができる.

本手法には, 下記のような 4 つの特徴がある. (1) 歌声を分離せず, ノイズ (伴奏音) が混在した状態をそのまま表現する. 人間は歌声を分離せずにそのまま認識できることを考えると, 人間の知覚の観点からも自然なやり方である. (2) 歌声と伴奏音の S/N 比を各フレームごとに推定可能なため, 伴奏音の変動に対して頑健である. (3) 単一の調波構造からスペクトル包絡を推定しないため, 高い F0 を持つ音に対しても頑健である. (4) F0 を持たない無声子音など, 他の音や音源に対しても, 調波フィルタを用いない歌声スペクトルテンプレートを用意することで容易に拡張できる.

3 定式化

本章では, 2 節で述べた手法の具体的な定式化について述べる.

3.1 確率的スペクトルテンプレート

歌声を含む混合音のスペクトル $y(f)$ は, 確率変数 Y_f から生成されると仮定する. ただし, f は対数軸での周波数を表し, s は対数軸でのスペクトルのパワーを表す. この確率変数 (の集合) Y_f を確率的スペクトルテンプレートと呼ぶ.

次に, Y_f は次式により 2 つの異なるスペクトルテンプレートに分割できると仮定する.

$Y_f = \log(\exp(Y_{v,f} + g_v) + \exp(Y_{n,f} + g_n))$ (1)
ただし, $Y_{v,f}$ は歌声のスペクトルを表し, 歌声スペクトルテンプレートと呼ばれ, $Y_{n,f}$ は歌声以外の音 (伴奏音) のスペクトルを表し, ノイズスペクトルテンプレートと呼ばれる. g_v と g_n はそれぞれのテンプレートの重みであり, それらを変化させることで歌声とその他の音の S/N 比を変化させることができる. なお, 式 (1) においては, 線形軸上でスペクトルの加

法性を仮定している.

$Y_{v,f}$ と $Y_{n,f}$ が, 次式のように, (対数周波数軸上で) 正規分布に従うと仮定する.

$$Y_{v,f} \sim \mathcal{N}(y; \mu_{v,f}, \sigma_{v,f}^2) \quad (2)$$

$$Y_{n,f} \sim \mathcal{N}(y; \mu_{n,f}, \sigma_{n,f}^2) \quad (3)$$

ここで, $\mathcal{N}(y; \mu, \sigma^2)$ は, 平均 μ , 分散 σ^2 の正規分布である. さらに, ソースフィルタモデルを仮定することで, 調波構造を持つ歌声 $Y_{v,f}$ は, 次式のように, 包絡の確率モデルと調波構造を表現するフィルタの対数軸上の加算で表現できると仮定する (図 2).

$$Y_{v,f} = Y'_{v,f} + \log H(f; f_0) \quad (4)$$

$$\sim \mathcal{N}(y; \mu'_{v,f} + \log H(f; f_0), \sigma_{v,f}^2) \quad (5)$$

$$H(f; f_0) = \sum_h \mathcal{N}(f; \log f_0 + \log h, \sigma_H^2) \quad (6)$$

ここで, $Y'_{v,f} \sim \mathcal{N}(y; \mu'_{v,f}, \sigma_{v,f}^2)$ は歌声のスペクトル包絡を表現する確率変数であり, 歌声包絡テンプレートと呼ぶ. また, $H(f; f_0)$ は F0 の値が f_0 のフィルタを表現し, 調波フィルタと呼ぶ. なお, 調波フィルタ $H(f; f_0)$ は確率変数ではないことに注意が必要である.

3.2 スペクトルテンプレートの加算の近似

式 (1) で表される確率的スペクトルテンプレート Y_f は, 解析的に計算することは困難であるので, 正規分布を用いて近似計算する. 関数 $l(x_1, x_2) = \log(\exp(x_1) + \exp(x_2))$ の $(x_1, x_2) = (\mu_{v,f} + g_v, \mu_{n,f} + g_n)$ における 2 次のテーラー展開を計算すると, パラメータ g_v, g_n, f_0 が固定された場合, x_1 と x_2 の重み付き加算で表現される. そのため,

$$Y_f = l(Y_{v,f}, Y_{n,f}) = \log(\exp(Y_{v,f}) + \exp(Y_{n,f})) \quad (7)$$

$$\text{は,} \quad Y_f \sim \mathcal{N}(y; \mu_f, \sigma_f^2) \quad (8)$$

$$\mu_f = \log(\exp(\mu_{v,f} + g_v) + \exp(\mu_{n,f} + g_n)) \quad (9)$$

$$\sigma_f^2 = \frac{(\exp(\mu_{v,f} + g_v))^2 \sigma_{v,f}^2 + (\exp(\mu_{n,f} + g_n))^2 \sigma_{n,f}^2}{(\exp(\mu_{v,f} + g_v) + \exp(\mu_{n,f} + g_n))^2} \quad (10)$$

のように近似することが出来る.

3.3 音素と F0 の推定

このモデルを使って音素と F0 を認識するためには, まず, それぞれの音素 i を表現する歌声包絡テンプレート θ_v^i とノイズスペクトルテンプレート θ_n を準備する必要がある. 観測スペクトル $y(f)$ が与えられたとき, 次式により $y(f)$ に含まれる音素 \hat{i} と F0 \hat{F}_0 を推定することができる.

$$(\hat{i}, \hat{F}_0) = \operatorname{argmax}_{i, f_0} \max_{g_v, g_n} \int_f \log \mathcal{N}(y(f); u_f, \sigma_f^2) df \quad (11)$$

ただし, u_f と σ_f^2 は, それぞれ式 (9) と (10) で定義される. また, 本稿の対象外ではあるが, 歌手名推定ができるように拡張したい場合は, 各歌手ごとに歌声包絡テンプレートを用意することで実現できる.

3.4 準ニュートン法によるパラメータ最適化

式 (11) を計算するためのパラメータ $\theta = (g_v, g_n, f_0)$ の最適化には, BFGS (Broyden-Fletcher-Goldfarb-

Shanno) 公式に基づく準ニュートン法を使用する。準ニュートン法は山登り法の一様であり、反復的にパラメータを更新する。本モデルにおいて、最小化すべき目的関数 $Q(\theta)$ は、

$$Q(\theta) = - \int_f \log \mathcal{N}(y(f); u_f, \sigma_f^2) df \quad (12)$$

で表される。ただし、 $y(f)$ は観測スペクトルである。

4 歌声包絡テンプレートの推定

式 (4) 中の歌声包絡テンプレート $Y'_{v,f}$ とノイズスペクトルテンプレート $Y_{n,f}$ は、学習データから推定する。一般に、調波構造を持つ歌声のスペクトルは、真のスペクトル包絡に対して、基本周波数の整数倍の周波数成分の点をサンプリングしたものと考えることができる。そのため、観測された歌声のスペクトル（調波構造）と、その元となるスペクトル包絡は一対多の関係になり得るので、単一フレームの調波構造から真のスペクトル包絡を推定することは困難である。本研究では、異なる F0 の値を持つ複数フレームの調波構造を用いることで、信頼性の高いスペクトル包絡を推定する。また、スペクトル包絡を一意に定めるのではなく、確率分布として推定するので、歌声の変動や学習データとテストデータの違いに対して頑健となる。

複数の調波構造からその元となるスペクトル包絡を推定する場合、フレームごとの音量の違いを考慮に入れる必要がある。そのため、本研究では各フレームの音量を正規化するためのパラメータを導入し、それも未知パラメータとして推定することでこの問題を解決する。

4.1 混合回帰分布

スペクトルテンプレートを表現するモデルとして、各回帰要素として線形回帰を使用した混合回帰モデル [10] を導入する。前章で述べたように、本手法においてはスペクトルテンプレートはある周波数 f における対数パワーの分布が正規分布で表現されるモデルを用いて定義される必要があるが、このモデルはその要件を満たしている。混合回帰モデルでは、スペクトルテンプレートの平均 $\mu_{v,f}$ と分散 $\sigma_{v,f}^2$ を

$$\mu_{v,f} = \sum_m G_m(f; \psi_m, \mu_m, \sigma_m^2) (a_m f + b_m) \quad (13)$$

$$\sigma_{v,f}^2 = \sum_m G_m(f; \psi_m, \mu_m, \sigma_m^2)^2 \beta_m^2 \quad (14)$$

として表現する。ただし、 $G_m(f; \psi_m, \mu_m, \sigma_m^2)$ はゲート関数の出力で、次式で定義される正規化ガウス関数 [11] を用いた。

$$G_m(f; \psi_m, \mu_m, \sigma_m^2) = \frac{\psi_m \mathcal{N}(f; \mu_m, \sigma_m^2)}{\sum_{m'} \psi_{m'} \mathcal{N}(f; \mu_{m'}, \sigma_{m'}^2)} \quad (15)$$

このモデルにおいて、未知パラメータは $\{\psi_m, \mu_m, \sigma_m^2, a_m, b_m, \beta_m^2\}$ であり、EM (Expectation and Maximization) 法により推定することが可能である。ただし、 ψ_m は、 $\psi_m \geq 0$ かつ $\sum_m \psi_m = 1$ である。

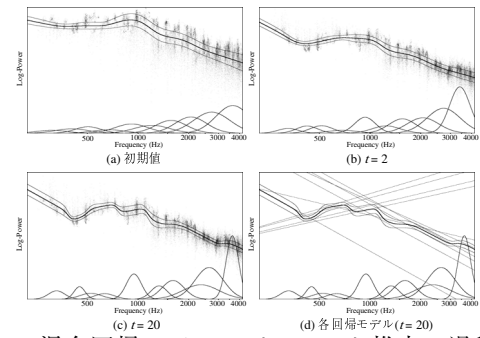


Fig. 3. 混合回帰モデルのパラメータ推定の過程の一例。各図の中心の太い線は混合回帰モデルの平均を表し、その上下の細い2本の線は標準偏差を表す。背景の細かい点は学習データの調波成分を表し、各図の下部の複数の山は、ゲート関数 $G_m(f; \psi_m, \mu_m, \sigma_m^2)$ を表す。

4.2 パラメータ推定

学習データとして与えられた I フレーム分の調波構造 $s_i (i = 1, \dots, I)$ の h 次倍音の周波数 $f_{i,h}$ とその対数パワー $y_{i,h}$ が、

$$s_n = \{(f_{i,1}, y_{i,1}), \dots, (f_{i,h}, y_{i,h}), \dots, (f_{i,H_i}, y_{i,H_i})\} \quad (16)$$

として表されるとする。この時、最大化したい尤度関数は、次式で表される。

$$L = \sum_i \sum_h \mathcal{N}(y_{i,h} + k_i; \mu_{v,f_{i,h}}, \sigma_{v,f_{i,h}}^2) \quad (17)$$

ここで、 k_i は各調波構造の音量を正規化するオフセットパラメータである。混合回帰モデルのパラメータと k_i を同時に最適化することは困難なので、それらを反復的に更新していく。図3は、パラメータ更新の過程を示す。ノイズスペクトルテンプレートについては、 $s_i (i = 1, \dots, I)$ を調波構造でなくスペクトルそのものと考え、同様に推定できる。

5 評価実験

本章では、提案法の性能を確認するために行った F0 と音素の同時推定の実験について述べる。

5.1 F0 と音素の同時推定

実験には、「RWC 研究用音楽データベース：ポピュラー音楽」[12] から選んだ10曲（男声3歌手、女声3歌手からなる）を用いた。音素推定の対象となる音素は日本語の5母音（/a/, /i/, /u/, /e/, /o/）とした。評価は、歌手ごとの6 fold cross validation により行った。各楽曲に対して、音素と F0 を手作業でアノテーションしたデータを学習用ラベル及び正解ラベルとして用いた。音素、F0 共に、対象の5母音を含むフレーム数に対する正しく認識できたフレーム数の割合を正解率として評価した。実験の際には、性別依存モデルを用いた。つまり、男声楽曲と女声楽曲で別々にテンプレートの集合（テンプレートモデル²）を学習し、識別の際には、男声テンプレートモデルと

²推定対象の複数の音素に対応する歌声包絡テンプレートと、ノイズスペクトルテンプレートの集合を、テンプレートモデルと呼ぶ。

Table 1. 音素と F0 の同時推定の実験結果 (正解率 [%]): 提案法の結果における↑は、比較法より性能が向上した場合を表す。

楽曲 *	性別	歌手	比較法		提案法	
			音素認識	F0 推定	音素認識	F0 推定
No. 4	男	A	31.1**	62.6**	73.5↑	58.9
No. 11	男	A	56.5	65.6	57.6↑	71.5↑
No. 9	男	B	47.5	65.5	43.4	43.3
No. 12	男	B	62.8	76.8	63.9↑	77.6↑
No. 6	男	C	51.5	69.2	60.4↑	80.8↑
No. 2	女	D	69.5	71.6	68.5	86.3↑
No. 16	女	D	62.7	78.2	65.4↑	82.6↑
No. 7	女	E	60.0	73.8	67.2↑	82.7↑
No. 18	女	E	64.1	73.5	70.2↑	87.6↑
No. 14	女	F	44.1	79.1	42.3	82.0↑
平均			55.0	71.6	61.2↑	75.3↑

* RWC 研究用音楽データベース: ポピュラー音楽 [12] の楽曲番号
 ** 異なる性別のモデルを誤って選択した楽曲

女声テンプレートモデルの両方で尤度を計算し、尤度が高いテンプレートモデルの結果を採用した。

比較法として、F0 に関しては PreFEst[13]³ を、音素推定に関しては文献 [1] の手法に基づいて分離した歌声から推定された MFCC を GMM により識別する手法を用いた。提案法では周波数解析に連続ウェーブレット変換を用いた。また、混合回帰モデルの混合数は 10 とした。比較法では、周波数解析に短時間フーリエ変換を用いた。また、MFCC の次元数は 12 次元とし、GMM の混合数は 32 とした。テンプレートの学習の際には、まず学習データとして使用する各楽曲に対して、歌声のみの音響信号と、歌声以外の伴奏音の音響信号 (カラオケトラック) を準備した。次に、各楽曲から、各音素に対して一つの歌声包絡テンプレートと、一つのノイズスペクトルテンプレートを学習した。識別の際に、各音素に対して尤度を計算する際は、その音素に対応するすべての歌声包絡テンプレートと、すべてのノイズスペクトルテンプレートの組み合わせに対して尤度を計算し、最も尤度の高い組み合わせの尤度を採用した。

実験結果を表 1 に示す。提案法により、10 曲の平均で音素推定は 6.2 ポイント、F0 推定は 3.7 ポイント性能が向上していることがわかる。音素推定では、10 曲中 7 曲で比較法より性能が向上している。特に No. 4 の楽曲では比較法では女声モデルの方が男声モデルより尤度が高くなっていたため、誤って女声モデルが使われてしまっているが、提案法では正しく男声モデルを選択できたので尤度が大幅に向上している⁴。F0 推定に関しては、10 曲中 8 曲で比較法より性能が向上している。一方で、No. 9 の楽曲は、提案法で F0 推定の正解率が 22.2 ポイントと大幅に低下している。この楽曲では、伴奏に使われているギターが大音量で鳴っており、そのギターの F0 を誤って推定してしまう場合が多かった。

³提案法では時間的連続性を考慮した処理を行っていないため、PreFEst においても各フレームの F0 の候補時間的連続性を考慮した後処理は行わなかった。

⁴なお、比較法において性別非依存のモデルを使用した場合は、No. 4 以外の楽曲では性別依存モデルの場合より性能が低下し、10 曲の平均でも性別依存モデルより 1 ポイント低い正解率だった。

6 まとめ

本稿では、多重奏の楽曲中の歌声の音素と F0 を同時に推定する手法について述べた。本手法の特徴は、歌声がその他の伴奏音と混ざった状態のスペクトルを、分離せずそのまま認識することにある。これは、人間は音を分離せずとも認識できるというアイデア [13] に基づいている。混合音を認識するための従来のやり方の多くは、構成するそれぞれの音を分離し、その後分離した音を認識するというアプローチだった。本研究のアプローチは背景のノイズに関する情報も活用するため、従来よりも性能を向上させることができる。

本手法は、音声認識の研究分野で知られる HMM 合成法 [14] と共通点がある。それは、クリーン音声 (歌声) のモデルとノイズのモデルを合成し、雑音下音声 (歌声) のモデルを作成する点である。HMM 合成法では、合成は学習段階で行われるのであらかじめ用意しておいた S/N 比でしか合成できなかったが、提案法は各フレームで S/N 比の推定を行うのでノイズの変動に対してロバストになるという利点がある。

本研究の最終的な目標は、歌詞を自動的に認識するシステムを実現することである。今後は、その実現を目指して、本フレームワークを拡張していく予定である。例えば、本稿で扱った 5 母音のみでなく、無声子音も含めたすべての音素について有効性を確認していく予定である。本研究の一部は CREST の支援を受けた。

7 参考文献

- [1] H. Fujihara, M. Goto, J. Ogata, K. Komatani, T. Ogata and H. G. Okuno: Automatic synchronization between lyrics and music CD recordings based on Viterbi alignment of segregated vocal signals, *Proc. ISM*, pp. 257–264 (2006).
- [2] H. Fujihara and M. Goto: Three Techniques for Improving Automatic Synchronization between Music and Lyrics: Fricative Sound Detection, Filler Model, and Novel Feature Vectors for Vocal Activity Detection, *Proc. ICASSP*, pp. 69–72 (2008).
- [3] 藤原弘将, 後藤真孝, 奥乃 博: 歌声の統計的モデル化とビタビ探索を用いた多重奏中のボーカルパートに対する音高推定手法, *情報処理学会論文誌*, **49** (2008).
- [4] M. Gruhne, K. Schmidt and C. Dittmar: Phoneme recognition in popular music, *Proc. ISMIR*, pp. 369–370 (2007).
- [5] D. Iskandar, Y. Wang, M.-Y. Kan and H. Li: Syllabic Level Automatic Synchronization of Music Signals and Text Lyrics, *Proc. ACM Multimedia*, pp. 659–662 (2006).
- [6] C. H. Wong, W. M. Szeto and K. H. Wong: Automatic lyrics alignment for Cantonese popular music, *Multimedia Syst.*, **4-5**, 307–323 (2007).
- [7] M.-Y. Kan, Y. Wang, D. Iskandar, T. L. Nwe and A. Shenoy: LyricAlly: Automatic Synchronization of Textual Lyrics to Acoustic Music Signals, *IEEE Trans. Audio, Speech, and Language Process.*, **16**, 338–349 (2008).
- [8] M. Ryyänänen and A. Klapuri: Transcription of the Singing Melody in Polyphonic Music, *Proc. ISMIR 2006*, pp. 222–227 (2006).
- [9] C. Sutton, E. Vincent, M. D. Plumbley and J. P. Bello: Transcription of vocal melodies using voice characteristics and algorithm fusion, *Proceedings of the MIREX 2006* (2006).
- [10] R. J. Jacobs, M. Jordan, S. J. Nowlan and G. E. Hinton: Adaptive mixtures of local experts, *Neural Computation*, **3**, 79–87 (1991).
- [11] L. Xu, M. I. Jordan and G. E. Hinton: An alternative model for mixtures of experts, *Adv. Neural Inf. Process. Syst.*, **7**, 633–640 (1994).
- [12] M. Goto, H. Hashiguchi, T. Nishimura and R. Oka: RWC Music Database: Popular, Classical, and Jazz Music Databases, *Proc. ISMIR*, pp. 287–288 (2002).
- [13] M. Goto: A Real-Time Music-Scene-Description System: Predominant-F0 Estimation for Detecting Melody and Bass Lines in Real-World Audio Signals, *Spe. Comm.*, **43**, 311–329 (2004).
- [14] M. J. F. Gales and S. Young: An improved approach to the hidden Markov model decomposition of speech and noise, *Proc. ICASSP*, pp. 835–838 (1997).