

歌唱における表現意図を考慮した 歌声 F_0 生成過程とその統計的モデリング

大石 康 智^{†1} 亀岡 弘 和^{†1,†2}
持橋 大地^{†3} 柏野 邦 夫^{†1}

歌声 F_0 軌跡における楽譜逸脱成分 (F_0 動的変動成分) を楽譜情報と分離して抽出することを目的として, F_0 軌跡の新しい生成過程モデルを提案する. 従来, オーバーシュートのような音符の立ち上がりに関する F_0 動的変動成分は 2 次系を用いてモデル化されたが, ビブラートやポルタメントのような意図的表現による変動成分および微細変動成分はすべて雑音としてモデル化された. 提案する F_0 生成過程は楽譜に記載される音符の並びを表現するノート指令信号と歌唱者の音楽的な表現意図を表す表現指令信号によって 2 次系が駆動されるモデルであり, ノート成分と表現成分を出力する. これらの成分とガウス性白色雑音に従う微細変動成分との和によって F_0 軌跡を記述する. そのモデルパラメータを推定する逆問題の解法アルゴリズムを導出し, 評価実験では, 推定される表現成分に歌唱者の音楽的な表現意図が含まれるかどうかを客観的かつ主観的に評価して, 提案モデルの有効性について議論する.

Statistical Modeling of Singing Voice F_0 Contours by Considering Musical Expressive Intentions

YASUNORI OHISHI,^{†1} HIROKAZU KAMEOKA,^{†1,†2}
DAICHI MOCHIHASHI^{†3} and KUNIO KASHINO^{†1}

We present a novel statistical model of singing voice fundamental frequency (F_0) contours for characterizing both musical-note information and various dynamic components such as *overshoot* and *vibrato*. Previous studies have modeled the dynamics using a second-order linear system and estimated the model parameters, but so far with limited success due to lack of flexibility for modeling the varied dynamics. Therefore we introduce the process of generating F_0 contours based on musical note and expression command functions, which generate the temporal attack of the note such as *overshoot* and the intended expression such as *vibrato*, respectively. Then we formulate a discrete-time stochastic process version of this model and propose a powerful framework for the estimation of the model parameters. In our experiment, the proposed method successfully decomposes F_0 contours into musical note and dynamic components.

1. はじめに

楽曲のメロディを歌った歌声の音高 (F_0) 軌跡には, そのメロディを構成する楽譜の音符の並びだけでなく, 楽譜に記載されない, “楽譜から逸脱した” 動的変動成分が含まれる (図 1). これらは, 発声器官の物理的制約に起因する成分 (特に, オーバーシュートや微細変動成分¹⁾⁻³⁾) と歌唱者の意図的表現による成分 (特に, ビブラートやポルタメント⁴⁾⁻⁶⁾) からなると考えられ, 知覚的には, 前者は人間らしさ・自然性に関係し, 後者は巧拙感に関係することがわかってきている⁷⁾⁻⁹⁾. さらに, 後者は意図して意図通りに表現できた場合と, 意図通りに制御できなかった場合とに分かれ, 習熟度に関連すると考えられる. ただし, これらの逸脱成分を F_0 軌跡から特徴抽出し, 歌唱者ごとにその特性を精緻に学習することまではまだ十分に検討されていない.

本研究の目的は, このような物理的制約もしくは意図的表現による楽譜逸脱成分を F_0 軌跡から楽譜情報と分離して抽出し, 歌唱者ごとにどのようなパターンをもちうるのか, 各パターンが文脈 (楽譜の音符列) にどう依存するかを計算機に学習させることである. ここでは, うろ覚えの状態では歌った歌声ではなく, 楽曲のメロディまたはその楽譜を既知として, 歌唱者なりに表情付けして歌った歌声を対象とする. 本研究は, 歌唱者の歌い方や個性, 癖を学習することを目指しており, 歌唱力評価や歌唱者識別, そして現在盛んに研究される歌声合成や歌声変換¹⁰⁾⁻¹⁸⁾ への応用が期待できる. 例えば, ある歌声を別の歌唱者の歌い方に変換して合成することが可能となるだろう. ここでは事前に歌唱者の歌い方が学習されるため, どんなメロディにもその歌い方を転写できることを特長とする.

本研究では, 物理モデルに基づいて F_0 軌跡の生成過程を記述し, そこから楽譜逸脱成分の特徴抽出に取り組む. 従来, 線形 2 次系を利用して歌声の F_0 に含まれる動的変動成分を制御するモデルが提案された¹⁹⁾⁻²²⁾. これらの研究では, 日本語の話し声の F_0 パターンを表現する藤崎モデル²³⁾ が参考にされた. 藤崎モデルは, 臨界制動 2 次系のインパルス応答とステップ応答を利用して, 日本語の句頭から句末に向けて緩やかに下降するフリーズ成分と, 語句に対応して急激に上昇下降するアクセント成分を表現し, これらを重畳することで

†1 日本電信電話株式会社 NTT コミュニケーション科学基礎研究所
NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation
†2 東京大学情報理工学系研究科
Graduate School of Information Science and Technology, The University of Tokyo
†3 大学共同利用機関法人 情報・システム研究機構 統計数理研究所
Institute of Statistical Mathematics

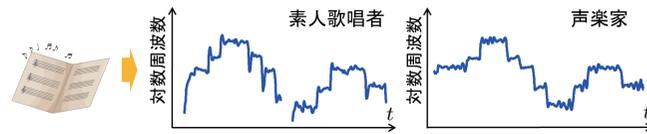


図1 歌声の基本周波数 (F_0) 軌跡とその動的変動成分

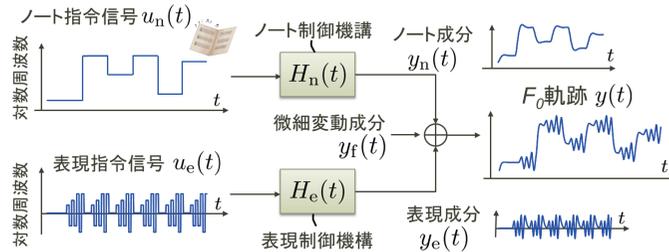


図2 提案する歌声の F_0 生成過程の概略図

F_0 を記述する．ただし，歌声のメロディに伴った急激な F_0 の上昇・下降の制御及び，ビブラートのような準周期的な振動は臨界制動系では表現できない．そのため，歌声の F_0 制御モデルでは2次系の伝達関数

$$G(s) = \frac{\Omega^2}{s^2 + 2\zeta\Omega s + \Omega^2} \quad (1)$$

における減衰率 ζ を調整することによって，指数減衰 ($\zeta > 1$)，減衰振動 ($0 < \zeta < 1$)，オーバーシュートを表現する)，臨界制動 ($\zeta = 1$)，定常振動 ($\zeta = 0$)，ビブラートを表現する) からなる様々な振動現象を表現する．文献 22) では，楽譜の音符列を表す階段状信号に式 (1) のインパルス応答を部分的に畳み込んで得られる F_0 軌跡を利用して，表情豊かな歌声合成音を実現した．しかしながら，制御パラメータ (減衰率 ζ と固有周波数 Ω) は手作業あるいは規則に基づいて決定された．

これに対し，我々は観測される F_0 軌跡から制御パラメータを推定する逆問題の解法を検討してきた．文献 24) は，楽譜の音符の並びを表す階段状信号を隠れマルコフモデル (HMM) からの出力で表現し，そこに2次系のインパルス応答が畳み込まれて，オーバーシュートのような動的変動成分が表現された．一方，ビブラートやブレバレーション，微細変動成分のような音符の音高が安定するときの楽譜逸脱成分はすべてガウス性白色雑音で表現された．そして，モデルパラメータの推定アルゴリズムを導出することで，楽譜逸脱成分の特徴抽出

を試みた．ただ，ビブラートや微細変動成分がすべてガウス性白色雑音としてモデル化されたため，歌唱者の意図的表現による動的変動成分 (ビブラートやポルタメントなど) を微細変動成分と分離して特徴付けられなかった．

本稿では，歌唱の意図的表現を特徴抽出するために，ノート指令信号と表現指令信号によって駆動される歌声 F_0 軌跡の生成過程を提案する (図 2)．ここで，ノート指令信号は楽譜に記載される音符の並びを表現する．一方，表現指令信号は歌唱者の音楽的な表現意図を矩形の細かい指令として表現する．ノート成分と表現成分はこれらの指令信号によって駆動されるノート制御機構と表現制御機構の出力である．ノート制御機構と表現制御機構はフィルタに相当し，2次系で表現される．ノート制御機構はオーバーシュートなどのノート (音符) の立ち上がり方を制御する．表現制御機構は矩形の細かい指令信号を制御してビブラートやポルタメントを生成する．ただし，表現制御機構は臨界制動系 ($\zeta = 1$ の場合) で構成される．微細変動成分は文献 22) にならって，10Hz 以上の不規則な振動成分を想定する．最終的に，対数スケールの F_0 軌跡 $y(t)$ (ここで， t は時間を表す) は，これら3つの成分の重ね合わせであると想定する．この F_0 生成過程を想定した理由は2点ある．

- 話声の F_0 生成過程を記述する藤崎モデルでは，甲状軟骨の二つの独立な運動 (平行移動と回転) に伴う声帯の長さの変化の合計が F_0 の時間的変化をもたらすと解釈され，平行移動運動に関係する成分をフレーズ成分，回転運動に関係する成分をアクセント成分とした²³⁾．これらの運動が歌声にも存在すると仮定し，フレーズ成分のような大域的な変化を表すノート成分と，アクセント成分のような歌唱者が意図的に制御できる表現成分から歌声 F_0 軌跡が構成されると考えた．
- 従来，ビブラートは正弦波で表現され，そのパラメータは Vibrato rate (ビブラートの速さ) と Vibrato extent (ビブラートの深さ) と考えられた^{11), 22)}．これらのパラメータが指数的に変化するモデルも提案された⁹⁾．しかし実際に，様々な歌唱者による多様な歌声の F_0 を観察した結果，これらのパラメータは時々刻々と変化するものの，指数的な変化であるとは言えない．この複雑な振動現象をモデル化して歌い方を認識するために，意図を表す矩形の細かい指令信号が臨界制動系によって制御されるビブラートの生成過程を想定した．

以降の節では，上記の歌声 F_0 軌跡の生成過程を離散時間表現して確率モデル化し，統計的手法に基づいてモデルパラメータの推定アルゴリズムを導出する．そして，提案モデルの有効性を客観的かつ主観的に評価する．文献 25) では，同じ枠組みで藤崎モデルを確率過程に基づいて統計モデル化する試みも行っているのだから参照されたい．

2. ノート指令と表現指令信号によって駆動される歌声 F_0 生成過程の確率モデル

歌声 F_0 軌跡の生成過程を離散時間表現し、確率過程に基づいて統計モデル化する。

2.1 F_0 生成過程の離散時間表現

連続時間領域で表現されるノート制御機構、表現制御機構の2次系の伝達関数の離散時間表現を得るために、従来法 [24] にならって後退差分変換を利用する。後退差分変換は、時間微分演算子 s を z 領域における後退差分演算子 $s \simeq (1 - z^{-1})/t_0$ に置き換える変換であり (t_0 は離散時間表現におけるサンプリング周期とする)、この変換によりノート制御機構の逆システムの伝達関数 $\mathcal{H}_n^{-1}(s)$ は z 領域で、

$$\mathcal{H}_n^{-1}(z) = a_2 z^{-2} + a_1 z^{-1} + a_0 \quad (2)$$

と書くことができる。ただし、

$$a_2 = \varphi^2, \quad a_1 = -2\varphi(\psi + \varphi), \quad a_0 = 1 + 2\varphi\psi + \varphi^2 \quad (3)$$

および、 φ, ψ は $\varphi = 1/(\Omega t_0)$, $\psi = \zeta$ と表現される。ここで、 k を離散時刻インデックスとし、ノート指令信号およびノート成分の離散時間表現をそれぞれ $u_n[k], y_n[k]$ とすると、 $y_n[k]$ は、ノート制御パラメータ φ, ψ によって特性が決まる拘束つき全極モデルからの出力

$$u_n[k] = a_0 y_n[k] + a_1 y_n[k-1] + a_2 y_n[k-2] \quad (4)$$

と見なすことができる。

同様に、表現指令信号 $u_e[k]$ と表現成分 $y_e[k]$ の関係も

$$u_e[k] = b_0 y_e[k] + b_1 y_e[k-1] + b_2 y_e[k-2] \quad (5)$$

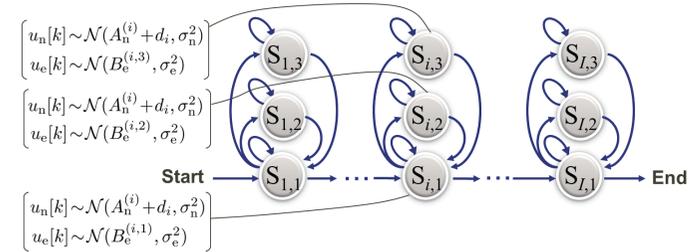
と書くことができる。ただし、 $b_2 = \xi^2$, $b_1 = -2\xi(1 + \xi)$, $b_0 = 1 + 2\xi + \xi^2$ であり、表現制御パラメータ ξ は $\xi = 1/(\Omega t_0)$ と表現される。微細変動成分 $y_f(t)$ の離散時間表現を $y_f[k]$ とすると、提案モデルによる歌声 F_0 軌跡の離散時間表現は、これら3つの成分の和

$$y[k] = y_n[k] + y_e[k] + y_f[k] \quad (6)$$

で与えられる。

2.2 F_0 生成過程の確率モデル化

ノート指令信号と表現指令信号はそれぞれ、楽譜に記載されるメロディの音符の並びと歌唱者の音楽的な表現意図を表す (図2)。これらの指令信号を表現するために、HMM を利



ノート指令信号 $\mu_n[k]$ と表現指令信号 $\mu_e[k]$ の例

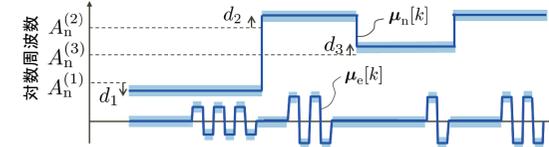


図3 隠れマルコフモデルによる指令信号の統計モデル化

用して、 $u_n[k]$ と $u_e[k]$ を確率モデル化する。 $o[k] := (u_n[k], u_e[k])^T$ を

$$o[k] \sim \mathcal{N}(\nu[k], \Upsilon), \quad \nu[k] := \begin{bmatrix} \mu_n[k] \\ \mu_e[k] \end{bmatrix}, \quad \Upsilon := \begin{bmatrix} \sigma_n^2 & 0 \\ 0 & \sigma_e^2 \end{bmatrix} \quad (7)$$

のように正規分布する確率変数と見なし、平均 $\nu[k]$ が図3のような状態遷移に伴って変化するモデルを考える。このように $o[k]$ を HMM でモデル化することにより、状態遷移の経路制限 (状態遷移確率の設定) を通して、 $\nu[k]$ に制約を与えることが可能となる。

具体的には、この HMM は $I \times J$ 個の状態集合 $S := \{S_{i,j}\}_{i=1,j=1}^{I,J}$ からなる。状態 $S_{i,j}$ では、 $\mu_n[k]$ は $A_n^{(i)} + d_i$ の値をとる。ここで、 $A_n^{(i)}$ は楽譜に記載されるメロディの i 番目の音符の音高を表し (本稿では、この値は楽譜から与えられるものとする)、 d_i はその音高からの推移 (音符の絶対音高からのずれであり、音高シフトパラメータと呼ぶ)、 I は歌唱するメロディに含まれる音符の総数を表す。図3に示す状態遷移により、 $\mu_n[k]$ は I 個の音符区間からなる階段状信号を表現する。一方、 $\mu_e[k]$ は $B_e^{(i,j)}$ の値をとり、これは歌唱者の表現意図を表すための、矩形の表現指令の大きさを表す。ここで、 i 番目の音符では、 $S_{i,1}$ を通らずして、状態 $S_{i,j}$ から別の状態 $S_{i,j'}$ ($j \neq j'$, $2 \leq j \leq J$, $2 \leq j' \leq J$) へ直接に遷移できない制約を設けることによって、 $\mu_e[k]$ は図3に示すような矩形信号を表現する。 J は表現指令を構成するための状態数を表す (図3は $J = 3$ である)。指令信号を生成する HMM

で $\mathbf{H} := [\mathbf{I}_K \ \mathbf{I}_K \ \mathbf{I}_K]$ とする．この場合，完全データの対数尤度は，

$$\log P(\mathbf{x}|\Theta) \stackrel{c}{=} \frac{1}{2} \log |\mathbf{\Lambda}^{-1}| - \frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{\Lambda}^{-1} (\mathbf{x} - \mathbf{m}) \quad (20)$$

$$\mathbf{m} := \begin{bmatrix} \mathbf{A}^{-1} \boldsymbol{\mu}_n \\ \mathbf{B}^{-1} \boldsymbol{\mu}_e \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{\Lambda}^{-1} := \begin{bmatrix} \mathbf{A}^T \boldsymbol{\Sigma}_n^{-1} \mathbf{A} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{B}^T \boldsymbol{\Sigma}_e^{-1} \mathbf{B} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \boldsymbol{\Sigma}_f^{-1} \end{bmatrix}$$

で与えられる．このとき，Q 関数 $Q(\Theta, \Theta')$ は，

$$Q(\Theta, \Theta') \stackrel{c}{=} \frac{1}{2} [\log |\mathbf{\Lambda}^{-1}| - \text{tr}(\mathbf{\Lambda}^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']) + 2\mathbf{m}^T \mathbf{\Lambda}^{-1} \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] - \mathbf{m}^T \mathbf{\Lambda}^{-1} \mathbf{m}] + \log P(\Theta) \quad (21)$$

となる．ここで， $\text{tr}(\cdot)$ は行列のトレースを表し， $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ と $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$ は条件付きガウス分布の性質より，

$$\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] = \mathbf{m} + \mathbf{\Lambda} \mathbf{H}^T (\mathbf{H} \mathbf{\Lambda} \mathbf{H}^T)^{-1} (\mathbf{y} - \mathbf{H} \mathbf{m}) \quad (22)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta'] = \mathbf{\Lambda} - \mathbf{\Lambda} \mathbf{H}^T (\mathbf{H} \mathbf{\Lambda} \mathbf{H}^T)^{-1} \mathbf{H} \mathbf{\Lambda} + \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']^T \quad (23)$$

と書ける．E ステップでは，直前のステップで更新されたパラメータを Θ' に代入し，上記に基づいて $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ と $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$ が算出される． \mathbf{y}_n ， \mathbf{y}_e ， \mathbf{y}_f に対応するように， $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ および $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$ を

$$\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] = \begin{bmatrix} \bar{\mathbf{x}}_n \\ \bar{\mathbf{x}}_e \\ \bar{\mathbf{x}}_f \end{bmatrix}, \quad \mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta'] = \begin{bmatrix} \mathbf{R}_n & * & * \\ * & \mathbf{R}_e & * \\ * & * & \mathbf{R}_f \end{bmatrix} \quad (24)$$

のように区分表現すると，Q 関数は

$$Q(\Theta, \Theta') \stackrel{c}{=} \frac{1}{2} \left[\log |\mathbf{A}^T \boldsymbol{\Sigma}_n^{-1} \mathbf{A}| + \log |\mathbf{B}^T \boldsymbol{\Sigma}_e^{-1} \mathbf{B}| + \log |\boldsymbol{\Sigma}_f^{-1}| - \text{tr}(\mathbf{A}^T \boldsymbol{\Sigma}_n^{-1} \mathbf{A} \mathbf{R}_n) + 2\boldsymbol{\mu}_n^T \boldsymbol{\Sigma}_n^{-1} \mathbf{A} \bar{\mathbf{x}}_n - \boldsymbol{\mu}_n^T \boldsymbol{\Sigma}_n^{-1} \boldsymbol{\mu}_n - \text{tr}(\mathbf{B}^T \boldsymbol{\Sigma}_e^{-1} \mathbf{B} \mathbf{R}_e) + 2\boldsymbol{\mu}_e^T \boldsymbol{\Sigma}_e^{-1} \mathbf{B} \bar{\mathbf{x}}_e - \boldsymbol{\mu}_e^T \boldsymbol{\Sigma}_e^{-1} \boldsymbol{\mu}_e - \text{tr}(\boldsymbol{\Sigma}_f^{-1} \mathbf{R}_f) \right] + \log P(\Theta) \quad (25)$$

と書き直せて，各パラメータについて M ステップの更新式を求めることができる．

1) 状態系列: Q 関数の中で $s := \{s_k\}_{k=1}^K$ に関する項は

$$\mathcal{I}_1(s) := -\frac{1}{2} \sum_{k=1}^K (\mathbf{o}[k] - \mathbf{c}_{s_k}[k])^T \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[k] - \mathbf{c}_{s_k}[k]) + \log P(s_1) + \sum_{k=2}^K \log P(s_k | s_{k-1})$$

となる．ここで， $\mathbf{o}[k] := ([\mathbf{A}\bar{\mathbf{x}}_n]_k, [\mathbf{B}\bar{\mathbf{x}}_e]_k)^T$ であり， $[\cdot]_k$ はベクトルの k 番目の要素を表す．これを最大化する状態系列 $\{s_k\}_{k=1}^K$ は動的計画法により効率的に解くことができる．まず，状態 $S_{1,1}$ について

$$\delta_1(S_{1,1}) = -\frac{1}{2} (\mathbf{o}[1] - \mathbf{c}_{S_{1,1}}[1])^T \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[1] - \mathbf{c}_{S_{1,1}}[1]) + \log P(S_{1,1}) \quad (26)$$

とおくと， $k = 2, 3, \dots, K$ について逐次的に $\delta_k(S_{i,j})$ を

$$\delta_k(S_{i,j}) = \max_{S_{i',j'}} \left[\delta_{k-1}(S_{i',j'}) - \frac{1}{2} (\mathbf{o}[k] - \mathbf{c}_{S_{i,j}}[k])^T \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[k] - \mathbf{c}_{S_{i,j}}[k]) + \phi_{S_{i',j'}, S_{i,j}} \right] \quad (27)$$

により計算できる．各ステップで選択される状態番号

$$\psi_k(S_{i,j}) = \operatorname{argmax}_{S_{i',j'}} [\delta_{k-1}(S_{i',j'}) + \phi_{S_{i',j'}, S_{i,j}}] \quad (28)$$

を記憶しておくことで， $k = K$ まで到着後に $s_{k-1} = \psi_k(s_k)$ ($k = K, K-1, \dots, 2$) により選択された状態番号を辿っていくと最適経路 s_1, \dots, s_K を得る．

2) ノート制御パラメータ: $\varphi^{(i)}$ と $\psi^{(i)}$ に関する事前分布を $\varphi^{(i)} \sim \mathcal{N}(\mu_\varphi, \sigma_\varphi^2)$ ， $\psi^{(i)} \sim \mathcal{N}(\mu_\psi, \sigma_\psi^2)$ とする．Q 関数の中で $\varphi^{(i)}$ と $\psi^{(i)}$ に関する項は，

$$\mathcal{I}_2(\varphi^{(i)}, \psi^{(i)}) = |\mathcal{T}_{S_{i,\cdot}}| \log(1 + 2\varphi^{(i)} \psi^{(i)} + (\varphi^{(i)})^2) - \frac{1}{2\sigma_n^2} \text{tr}((\mathbf{A}^{(i)})^T \mathbf{A}^{(i)} \mathbf{R}_n) + \frac{1}{\sigma_n^2} ([\boldsymbol{\mu}_n]_{\mathcal{T}_{S_{i,\cdot}}})^T \mathbf{A}^{(i)} \bar{\mathbf{x}}_n - \frac{1}{2\sigma_\varphi^2} (\varphi^{(i)} - \mu_\varphi)^2 - \frac{1}{2\sigma_\psi^2} (\psi^{(i)} - \mu_\psi)^2 \quad (29)$$

$$\mathcal{T}_{S_{i,\cdot}} = \{k | s_k \in \{S_{i,1}, \dots, S_{i,J}\}\}$$

となる．ここで， $|\mathcal{T}|$ は集合 \mathcal{T} の要素数を表す．また， $[\boldsymbol{\mu}]_{\mathcal{T}}$ は，集合 \mathcal{T} の要素を添え字として，その添え字に相当する $\boldsymbol{\mu}$ の要素を取り出した部分ベクトルを表す．今，

$$\mathbf{U}_2 := \begin{bmatrix} 1 & & & \mathbf{O} \\ -2 & 1 & & \\ 1 & -2 & 1 & \\ & \ddots & \ddots & \ddots \\ \mathbf{O} & & 1 & -2 & 1 \end{bmatrix}, \quad \mathbf{U}_1 := \begin{bmatrix} 2 & & & \mathbf{O} \\ -2 & 2 & & \\ 0 & -2 & 2 & \\ & \ddots & \ddots & \ddots \\ \mathbf{O} & & 0 & -2 & 2 \end{bmatrix}, \quad \mathbf{U}_0 := \begin{bmatrix} 1 & & & \mathbf{O} \\ 0 & 1 & & \\ 0 & 0 & 1 & \\ & \ddots & \ddots & \ddots \\ \mathbf{O} & & 0 & 0 & 1 \end{bmatrix}$$

として、式 (4) から、 $A^{(i)}$ は、

$$A^{(i)} = (\varphi^{(i)})^2 [U_2]_{\mathcal{I}_{S_i}} + \varphi^{(i)} \psi^{(i)} [U_1]_{\mathcal{I}_{S_i}} + [U_0]_{\mathcal{I}_{S_i}} \quad (30)$$

と表現される。ここで、 $[U]_{\mathcal{I}}$ は集合 \mathcal{I} の要素を添え字として、行列 U からその添え字に相当する行ベクトルを取り出して構成される部分行列を意味する。ニュートン・ラフソン法を利用して、 $\mathcal{I}_2(\varphi^{(i)}, \psi^{(i)})$ を最大化する $\varphi^{(i)}$ と $\psi^{(i)}$ が数値的に導出される。

3) 表現制御パラメータ: ξ に関する事前分布を $\xi \sim \mathcal{N}(\mu_\xi, \sigma_\xi^2)$ とする。Q 関数の中で ξ に関係する項は、

$$\mathcal{I}_3(\xi) = K \log(1 + 2\xi + \xi^2) - \frac{1}{2\sigma_e^2} \text{tr}(\mathbf{B}^T \mathbf{B} \mathbf{R}_e) + \frac{1}{\sigma_e^2} \mu_e^T \mathbf{B} \bar{\mathbf{x}}_e - \frac{1}{2\sigma_\xi^2} (\xi - \mu_\xi)^2 \quad (31)$$

となる。ニュートン・ラフソン法を利用して、 $\mathcal{I}_3(\xi)$ を最大化する ξ が数値的に導出される。

4) その他のパラメータ: d_i と $B_e^{(i,j)}$ に関して、それぞれ事前分布を $d_i \sim \mathcal{N}(0, \sigma_d^2)$ と $B_e^{(i,j)} \sim \mathcal{N}(\mu_{B^{(i,j)}}, \sigma_B^2)$ とする。残されたパラメータの更新式は

$$d_i = \frac{1}{|\mathcal{I}_{S_i}| + \sigma_n^2 / \sigma_d^2} \sum_{k \in \mathcal{I}_{S_i}} ([A \bar{\mathbf{x}}_n]_k - A_n^{(i)}) \quad (32)$$

$$B_e^{(i,j)} = \frac{1}{|\mathcal{I}_{S_i,j}| / \sigma_e^2 + 1 / \sigma_B^2} \left(\sum_{k \in \mathcal{I}_{S_i,j}} \frac{[B \bar{\mathbf{x}}_e]_k}{\sigma_e^2} + \frac{\mu_{B^{(i,j)}}}{\sigma_B^2} \right), \quad \mathcal{I}_{S_i,j} = \{k | s_k = S_{i,j}\} \quad (33)$$

$$\sigma_n^2 = (\text{tr}(\mathbf{A}^T \mathbf{A} \mathbf{R}_n) - 2\mu_n^T \mathbf{A} \bar{\mathbf{x}}_n + \mu_n^T \mu_n) / K \quad (34)$$

$$\sigma_e^2 = (\text{tr}(\mathbf{B}^T \mathbf{B} \mathbf{R}_e) - 2\mu_e^T \mathbf{B} \bar{\mathbf{x}}_e + \mu_e^T \mu_e) / K \quad (35)$$

$$\sigma_f^2 = \text{tr}(\mathbf{R}_f) / K \quad (36)$$

と導出される。

4. 実装方法

前節で導出したパラメータ推定アルゴリズムにおける実際の実装方法について述べる。

4.1 実測 F_0 軌跡からパラメータ推定方法

提案アルゴリズムを用いて、実測 F_0 軌跡からパラメータを推定する問題を考える。式 (18) では、全区間で F_0 が観測されていることが暗に想定されているが、実際には F_0 は有声音が発せられるときのみ観測可能であり、無声音の区間では F_0 は通常観測できない。したがって、実測 F_0 軌跡からのパラメータ推定を行うためには一般に欠損データの問題を扱う必要がある。この類の問題は不完全データ問題に他ならず、EM アルゴリズムにより扱う

ことができる。文献 25) を参照して、 $\mathbf{y} \in \mathbb{R}^K$ を実測 F_0 データと欠損 F_0 データからなる完全データとし、実測 F_0 データを並べたベクトルを $\mathbf{y}_{\text{obs}} \in \mathbb{R}^{K'}$ ($K' \leq K$) とする。 \mathbf{y} と \mathbf{y}_{obs} との関係は、各行に 1 が一個あり残りはすべて 0 であるような $K' \times K$ のバイナリ行列 M を用いて $\mathbf{y}_{\text{obs}} = M\mathbf{y}$ と表される。これを利用して EM の各ステップを導くことができる。導出は割愛するが、E ステップでは \mathbf{y} を

$$\mathbf{y} \leftarrow \boldsymbol{\mu} + \Sigma M^T (M \Sigma M^T)^{-1} (\mathbf{y}_{\text{obs}} - M\boldsymbol{\mu}) \quad (37)$$

により更新し、M ステップではその \mathbf{y} を用いて

$$\Theta \leftarrow \underset{\Theta}{\text{argmax}} \log P(\mathbf{y}|\Theta)P(\Theta) \quad (38)$$

を実行すれば良い。なお自明ではあるが、欠損区間がない場合を表す $M = I_K$ においては式 (37) は $\mathbf{y} \leftarrow \mathbf{y}_{\text{obs}}$ となり、 \mathbf{y}_{obs} をそのまま完全データと見なして良いという意味になる。

4.2 パラメータの更新方法

EM 法に基づいてパラメータ推定アルゴリズムを導出したものの、実際は推定すべきパラメータ数が多いため、各パラメータを反復更新しても、局所解に収束してしまう。これを抑制するために、本稿では表現制御パラメータ ξ を固定する。したがって、意図的表現に関する情報は表現指令信号のみを推定する。また、以下の順序で分散パラメータを小さな値に固定する制約の下で、各パラメータを反復して更新する。

- (1) σ_e^2, σ_f^2 を固定して (例えば、どちらも 1)、E ステップと M ステップを反復する。
- (2) σ_n^2, σ_f^2 を固定して、E ステップと M ステップを反復する。
- (3) σ_n^2, σ_e^2 を固定して、E ステップと M ステップを反復する。

5. 評価実験

提案モデルの有効性を、人工的に作成した F_0 軌跡、および実測の F_0 軌跡を用いて客観的かつ主観的に評価する。

5.1 人工的に作成した F_0 軌跡を用いた評価

まず、パラメータ推定アルゴリズムの挙動を調べるために、文献 22) の歌声 F_0 制御モデルを用いて、人工的に作成した F_0 軌跡をテストデータとして動作実験を行った。事前に童謡や歌謡曲から抜粋したメロディの MIDI (合計 16 種類、抜粋部分の長さは平均 11.4 秒) に基づいて、対数周波数 (単位は cent, 半音は 100cent に相当する) 軸上で階段状に変化する音符列を作成する (サンプリング周期は 5ms, 休符区間はその前の音符の音高

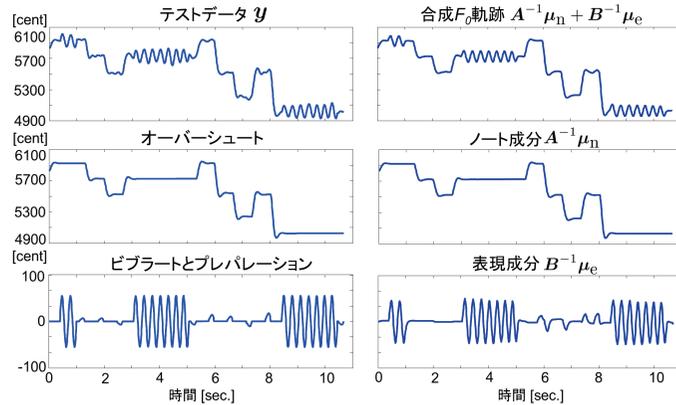


図4 人工的に作成した F_0 軌跡に対するパラメータ推定結果

を伸ばす)。これに、文献 22) で調査されたプロ歌唱者の F_0 軌跡を生成するためのパラメータを用いて、楽譜逸脱成分を付加する (図 4 の左側)。具体的には、オーバーシュートは $\zeta = 0.570$, $\Omega = 0.0363$ [rad/ms] とし、プレパレーションは $\zeta = 0.675$, $\Omega = 0.0308$ [rad/ms] とし、ビブラートは Vibrato extent = 83.0 cent, Vibrato rate = 6.25 Hz とし、微細変動成分は白色雑音をカットオフ周波数 10 Hz のローパスフィルタに通し、振幅の平均値を 20 cent に調整したものとし、キーシフトは +25cent, 0cent, -25cent の 3 パターンとした。以上より、合計 48 個のテストデータを得る。

一方、パラメータ推定アルゴリズムの設定値を以下に示す。前節の (1)~(3) における各反復回数は 1000 回とした。I と $\{A_n^{(i)}\}_{i=1}^I$ は抜粋したメロディの MIDI 情報から与えられる。表現指令信号を構成するための状態数 J は 5 とした。HMM の状態遷移確率は、 $\phi_{S_{i,1}, S_{i,1}} = \log(0.9999 \times (J - 1)/J)$, $\phi_{S_{i,1}, S_{i,j}} = \log(0.9999/J)$, $\phi_{S_{i,1}, S_{i+1,1}} = \log(0.0001)$, $\phi_{S_{i,j}, S_{i,j}} = \log(0.9999)$, $\phi_{S_{i,j}, S_{i,1}} = \log(0.0001)$, ($1 \leq i \leq I$, $2 \leq j \leq J$) とした。表現制御パラメータは $\xi = 3$ に固定した。パラメータの事前分布における固定パラメータは、 $\mu_\varphi = 6$, $\sigma_\varphi^2 = 0.1$, $\mu_\psi = 0.6$, $\sigma_\psi^2 = 0.02$, $\sigma_d^2 = 2500$, $\sigma_B^2 = 100$, $\mu_{B^{(i,1)}} = 0$, $\mu_{B^{(i,2)}} = 30$, $\mu_{B^{(i,3)}} = -30$, $\mu_{B^{(i,4)}} = 60$, $\mu_{B^{(i,5)}} = -60$, ($1 \leq i \leq I$) とした。これらは文献 22) と予備実験に基づいて決定した結果である。

図 4 はテストデータに対するパラメータの推定結果例を示す。左側がテストデータであり、右側が推定結果である。テストデータに近い軌跡が推定されたことを定性的に確認でき

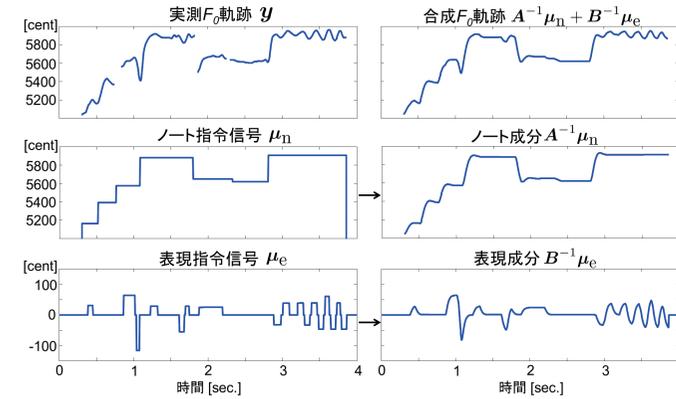


図5 実測 F_0 軌跡に対するパラメータ推定結果

る。二乗平均平方根 (Root Mean Square, RMS) を計算したところ、ノート成分とオーバーシュートを含む信号との RMS は 7.47cent, 表現成分とビブラートおよびプレパレーションを含む信号との RMS は 9.51cent であった。対数周波数上の 100cent は半音に相当し、RMS がその 10% 以下に収まっていることから、人工的な F_0 軌跡を用いた場合、パラメータ推定アルゴリズムの局所解問題を抑制できていることを確認できた。

5.2 実測 F_0 軌跡を用いた評価

次に、実測の F_0 軌跡に対する推定結果を図 5 に示す。実測 F_0 軌跡として、「RWC 研究用音楽データベース：ポピュラー音楽」(RWC-MDB-P-2001)²⁶⁾ における、歌手名：緒方智美、曲番号：No.07、曲名：PROLOGUE のメロディの F_0 を手作業でラベル付した結果²⁷⁾ を用いた。本来ならばこの楽曲の音響信号から F_0 を推定すべきであるが、今回は提案手法の性能の上限を調べるためにこのようなデータを利用した。 F_0 は 10ms ごとにラベル付されているので、アップサンプリングによって、5ms ごとの実測 F_0 軌跡を得た。同楽曲の楽譜情報は「RWC 研究用音楽データベース：ポピュラー音楽」の MIDI データを利用した。図 5 は推定結果の一部であり、実測 F_0 軌跡からノート指令信号と表現指令信号が推定され、それらの制御機構のインパルス応答が畳み込まれて、ノート成分と表現成分が得られる。これらを足し合わせたものを合成 F_0 軌跡とし、定性的には実測 F_0 軌跡に近い軌跡が得られたことがわかる。また、微細変動成分パラメータである σ_ψ^2 の値は 187.7 であるため、振幅が 13cent 程度の微細な変動成分が推定されており、先行研究の知見³⁾ と整合する。特徴的

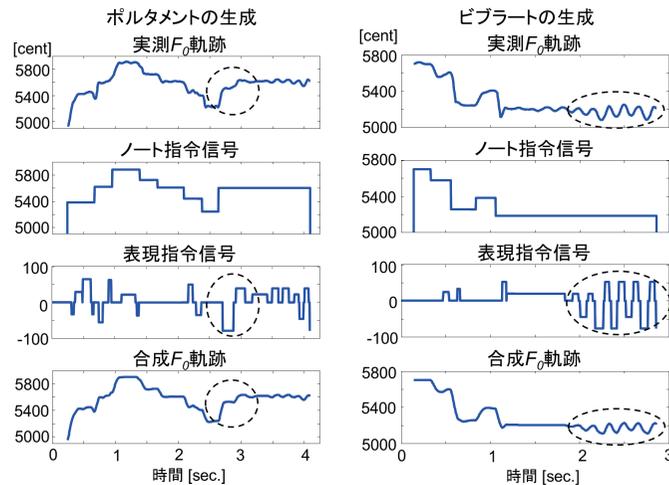


図6 ポルタメントとビブラートが生成される様子

なことは、4節で述べたように、無声音の区間を欠損データとして扱っているが、合成 F_0 軌跡ではその部分が周囲の F_0 から補間されていることである。

具体的に、推定されたノート指令信号と表現指令信号を見てみよう。図6の左図はポルタメントが生成される様子を示す。楽譜(MIDI)では一つの音を伸ばし続けるよう記述されているものの、歌手の表現意図によって、半音だけ滑らかに変化させる様子が実測 F_0 軌跡に観測される。このとき提案法では、ノート指令信号は同一音符区間として推定され、表現指令信号は点線枠内に示されるような音高を変化させる指令が推定される。一方、図6の右図はビブラートが生成される様子を示す。実測 F_0 軌跡の最後の音符区間にビブラートが観測される。ビブラートが必ずしも正弦波に従う挙動ではないことがわかる。このとき、表現指令信号として、部分的に矩形化する指令が推定された。ポルタメントとビブラートの生成のどちらにおいても合成 F_0 軌跡は実測 F_0 軌跡に近い軌跡が得られている。

5.3 表現指令パターンの学習

前節より、音符区間におけるポルタメントやビブラートのような表現意図を表現指令信号として推定できることがわかった。そうすると、この歌唱者の表現指令信号に頻出する標準パターンを学習したくなる。楽曲「PROLOGUE」には同一音高、同一音長の音符が複数存在するため、ある同一音符区間における表現指令入力信号を描画した結果を図7に示す。こ

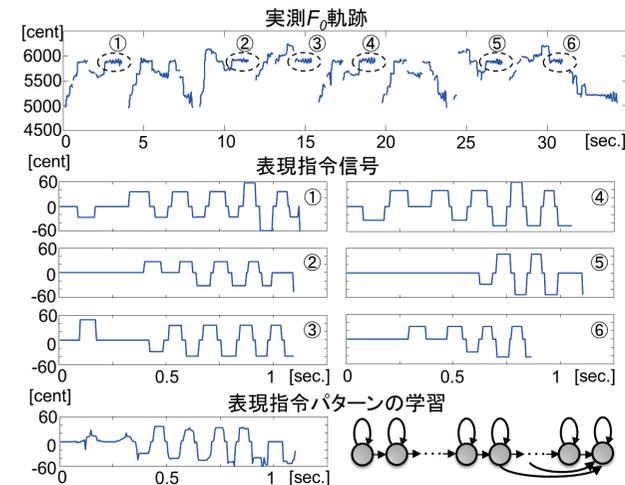


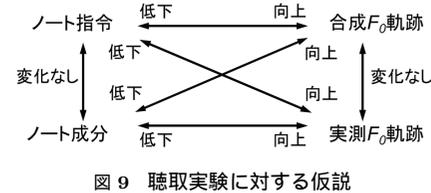
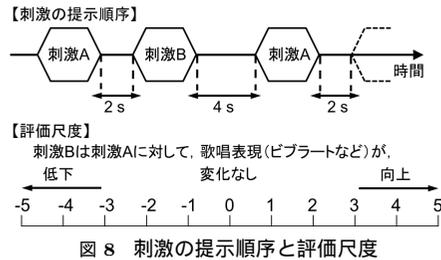
図7 推定された表現指令信号を用いた表現指令パターンの学習

こでは6つの表現指令信号を示した(楽曲全体で18個ある)。この音符区間ではビブラートをかけるため、似たような表現指令が得られる。これらの表現指令信号の標準パターンを得るために、図7の下部に示す Left-to-right 型の HMM を利用して、指令信号の軌跡の学習を試みた。この音符区間の表現指令信号の平均長に基づいて、状態数は220とした。ただ、平均長よりも短い表現指令信号もあるため、最後の50状態は最終状態への遷移確率を持つように HMM のトポロジを構成した。出力確率分布は単一ガウス分布として HMM のパラメータを推定し、ガウス分布の平均値を描画した結果を図7の下部に示す。0.5秒頃からビブラートの振動が開始し、その振幅が時間とともに右下がりに変動するパターンが学習された。今後の課題は、このように推定された表現指令信号から歌唱者の表現意図を音符区間ごとに大規模に学習することである。

5.4 聴取実験に基づく提案モデルの評価

最後に、提案モデルとそのパラメータ推定アルゴリズムによる歌唱の意図的表現の抽出について、主観的な聴取実験に基づいて評価する。5.2節で示した実測 F_0 軌跡から推定される、4種類の信号軌跡を準備する。図5と対応付けて見てほしい。

- (1) ノート指令: μ_n
- (2) ノート成分: $A^{-1}\mu_n$



- (3) 合成 F_0 軌跡: $A^{-1}\mu_n + B^{-1}\mu_e$
 (4) 実測 F_0 軌跡: y

これらの信号軌跡と YAMAHA の歌声合成ソフトウェア Vocaloid3²⁸⁾ (歌手は Vocaloid2 用の VY1)^{*1} を利用して、歌声を合成する。具体的には、推定されたノート指令信号に基づいて、発音時刻ごとにノートと歌詞を入力する (今回はすべての音符を「な、/na/」で発音する)。次に、ピッチバンドセンシティビティを 12 (1 オクターブ) に設定して、(1) ~ (4) の信号軌跡それぞれに対してピッチバンドの値を計算し、入力する。音量はすべて一定とする。以上の手順によって、4 種類の歌声が合成される。最後にこれらの音響信号を 4 小節ごとに切り出した全 96 サンプル (歌声ごとに 24 サンプルが切り出される) を刺激として用いる。図 8 に示すように、ランダムに刺激を選び出して、刺激 A、刺激 B と提示し、これを 50 セット準備する。被験者数は成人男性 5 名とし、各被験者は、刺激 B は刺激 A に対して歌唱表現が向上したか否か (ビブラートやポルタメントを付加し、滑らかに歌っているか否か) を 11 段階で評価する。この聴取実験の結果に対する仮説を図 9 に示す。着目すべきは、

- 合成 F_0 軌跡による歌声が、実測 F_0 軌跡による歌声と比べて、歌唱表現に変化がなく評価されること
- 合成 F_0 軌跡による歌声が、ノート指令による歌声およびノート成分による歌声と比べて、歌唱表現が向上したと評価されること

と仮説を立てている点である。

評価結果を表 1 に示す。被験者ごとに評価値の分散が異なるので、各被験者で値を正規化し、それを 5 人の被験者で平均した結果が表 1 である。仮説のとおり、刺激

表 1 聴取実験結果：評価値を被験者ごとに正規化した後の平均値

刺激 A \ 刺激 B	ノート指令	ノート成分	合成 F_0 軌跡	実測 F_0 軌跡
ノート指令		0.151	0.479	0.434
ノート成分	-0.106		0.560	0.632
合成 F_0 軌跡	-0.621	-0.524		0.123
実測 F_0 軌跡	-0.679	-0.540	0.091	

A { ノート指令, ノート成分 } → 刺激 B { 合成 F_0 軌跡, 実測 F_0 軌跡 } のとき評価値は正となり、刺激 A { 合成 F_0 軌跡, 実測 F_0 軌跡 } → 刺激 B { ノート指令, ノート成分 } のとき評価値は負となり、合成音における歌唱表現の有無を評価できていることがわかる。なおかつ、刺激 A { 合成 F_0 軌跡 } → 刺激 B { 実測 F_0 軌跡 } および刺激 A { 実測 F_0 軌跡 } → 刺激 B { 合成 F_0 軌跡 } のとき、評価値が他に比べて 0 に近い値となった。このことから、実測 F_0 軌跡に含まれる歌唱の意図的表現が、提案法の表現成分として推定されることを主観的な観点から確認できた。ただし、合成 F_0 軌跡による歌声と実測 F_0 軌跡による歌声を比較したときの評価値が 0 でないことから、残された微細変動成分について今後詳細に分析する必要がある。また本実験では対象曲が 1 曲なので、今後は様々な歌唱者が多様な楽曲を歌った歌声を対象として、提案モデルの有効性を大規模に評価する必要がある。

6. おわりに

本稿では、歌声 F_0 軌跡に含まれる物理的制約もしくは意図的表現による楽譜逸脱成分 (F_0 動的変動成分) を楽譜情報と分離して抽出することを目的として、ノート指令信号と表現指令信号によって駆動される F_0 軌跡の生成過程を記述し、そのモデルパラメータを推定する逆問題の解法を提案した。提案する F_0 生成過程は楽譜に記載される音符の並びを表現するノート指令信号と歌唱者の音楽的な表現意図を表す表現指令信号が 2 次系によって制御され、それぞれノート成分と表現成分を出力する。これらの成分とガウス性白色雑音に従う微細変動成分との和によって F_0 軌跡が記述される。評価実験では、導出したパラメータ推定アルゴリズムによって推定される表現成分に、歌唱者の意図的表現が含まれることを客観的かつ主観的に確認し、提案モデルの有効性を示した。今後の課題は、抽出された楽譜

*1 <http://www.vocaloid.com/lineup/>

逸脱成分が歌唱者ごとにどのようなパターンをもちうるのか、各パターンが文脈（楽譜の音符列）にどう依存するかを学習することである。また、大規模な聴取実験を行うことで、常に人間の歌声知覚も考慮しながら研究を進める必要がある。さらには、 F_0 （音高）だけでなく、音韻や音量における楽譜逸脱成分についても検討を進めたい。

謝辞 本研究に対し、有益なご助言を頂いた中野允裕氏（NTT CS 研）に感謝致します。

参 考 文 献

- 1) Sundberg, J.: *The Science of the Singing*, Northern Illinois University Press (1987).
- 2) de Krom, G. and Bloothoof, G.: Timing and Accuracy of Fundamental Frequency Changes in Singing, *In Proc. ICPHS95*, pp.206–209 (1995).
- 3) Akagi, M. and Kitakaze, H.: Perception of Synthesized Singing Voices with Fine Fluctuations in Their Fundamental Frequency Contours, *in Proc. ICSLP 2000*, pp. 458–461 (2000).
- 4) Seashore, C. E.: A Musical Ornament, the Vibrato, *In Psychology of Music*, McGraw-Hill Book Company, pp.33–52 (1938).
- 5) Kenji Kojima, M. Y. and Nakayama, I.: Variability of Vibrato -A Comparative Study between Japanese Traditional Singing and Bel Canto-, *In Proc. Speech Prosody 2004*, pp.151–154 (2004).
- 6) Nakayama, I.: Comparative Studies on Vocal Expressions in Japanese Traditional and Western Classical- Style Singing, Using a Common Verse, *In Proc. ICA 2004*, pp.1295–1296 (2004).
- 7) Saitou, T. et al.: Speech-To-Singing Synthesis: Converting Speaking Voices to Singing Voices by Controlling Acoustic Features Unique to Singing Voices, *in Proc. WASSPA 2007*, pp.215–218 (2007).
- 8) Saitou, T. et al.: Acoustic and Perceptual Effects of Vocal training in Amateur Male Singing, *in Proc. EUROSPEECH 2009*, pp.832–835 (2009).
- 9) 右田尚人ほか：ヴィブラート歌唱における基本周波数制御に有効な特徴量の検討，音講論集，3-P-28，pp.395–396 (2010).
- 10) 河原英紀，片寄晴弘：高品質音声分析変換合成システム STRAIGHT を用いたスクラップ生成研究の提案，情報処理学会論文誌，Vol.43, No.2, pp.208–218 (2002).
- 11) Nakano, T. et al.: An Automatic Singing Skill Evaluation Method for Unknown Melodies Using Pitch Interval Accuracy and Vibrato Features, *Proc. ICSLP 2006*, pp.1706–1709 (2006).
- 12) Mayor, O. et al.: The Singing Tutor: Expression Categorization and Segmentation of the Singing Voice, *in Proc. AES 121st Convention* (2006).
- 13) 中野倫靖，後藤真孝：VocaListener：ユーザ歌唱の音高および音量を真似る歌声合成システム，情報処理学会論文誌，Vol.52, No.12, pp.3853–3867 (2011).
- 14) Bonada, J. et al.: Synthesis of the Singing Voice by Performance Sampling and Spectral Models, *IEEE Signal Processing Magazine*, Vol.24, pp.67–79 (2007).
- 15) Kako, T. et al.: Automatic Identification for Singing Style Based on Sung Melodic Contour Characterized in Phase Plane, *in Proc. ISMIR 2009*, pp.393–397 (2009).
- 16) Fukayama, S. et al.: Orpheus: Automatic Composition System Considering Prosody of Japanese Lyrics, *in Proc. ICEC 2009*, pp.309–310 (2009).
- 17) Nakano, T. et al.: VocaListener2: A Singing Synthesis System Able to Mimic a User's Singing in terms of Voice Timbre Changes as Well as Pitch and Dynamics, *in Proc. ICASSP 2011*, pp.453–456.
- 18) Mase, A. et al.: HMM-based singing voice synthesis system using pitch-shifted pseudo training data, *in Proc. INTERSPEECH 2010*, pp.845–848.
- 19) 柏野邦夫ほか：パート譜を用いたボーカル音分離システム，音講論集，2-9-1，pp.625–626 (1998).
- 20) Mori, H., Odagiri, W. and Kasuya, H.: F_0 dynamics in singing: Evidence from the data of a baritone singer, *IEICE Trans. Inf. and Syst.*, Vol.E87-D, No.5, pp. 1086–1092 (2004).
- 21) Minematsu, N. et al.: Prosodic Modeling of Nagauta Singing and Its Evaluation, *Proc. SpeechProsody 2004*, pp.487–490 (2004).
- 22) Saitou, T., Unoki, M. and Akagi, M.: Development of an F_0 control model based on F_0 dynamic characteristics for singing-voice synthesis, *Speech Communication*, Vol.46, pp.405–417 (2005).
- 23) Fujisaki, H.: A note on the physiological and physical basis for the phrase and accent components in the voice fundamental frequency contour, *Vocal Physiology: Voice Production, Mechanisms and Functions*, (O. Fujimura, ed.), Raven Press, pp. 347–355 (1988).
- 24) 大石康智ほか：歌声 F_0 系列からの楽譜逸脱成分の抽出 - 動特性モデルに基づく楽譜との時間的対応付け - ，音講論集，1-8-19，pp.279–282 (2011).
- 25) Kameoka, H. et al.: A statistical model of speech F_0 contours, *in Proc. SAPA 2010*, pp.43–48.
- 26) 後藤真孝，橋口博樹，西村拓一，岡 隆一：RWC 研究用音楽データベース：研究目的で利用可能な著作権処理済み楽曲・楽器音データベース，情報処理学会論文誌，Vol.45, No.3, pp.728–738 (2004).
- 27) Goto, M.: AIST Annotation for the RWC Music Database, *in Proc. ISMIR 2006* (2006).
- 28) Kenmochi, H. and Ohshita, H.: VOCALOID - Commercial Singing Synthesizer based on Sample Concatenation, *in Proc. INTERSPEECH 2007*, pp.4010–4011 (2007).