

[ポスター講演] 調とリズムを考慮した 階層隠れセミマルコフモデルに基づく歌声の自動採譜

錦見 亮[†] 中村 栄太[†] 後藤 真孝^{††} 糸山 克寿[†] 吉井 和佳^{†,††}

[†] 京都大学大学院情報学研究所

^{††} 産業技術総合研究所 (AIST)

^{†††} 理化学研究所 革新知能統合研究センター (AIP)

E-mail: [†]{nishikimi,enakamura,itoyama,yoshii}@sap.ist.i.kyoto-u.ac.jp, ^{††}m.goto@aist.go.jp

あらまし 本稿では歌声 F0 軌跡から音楽的に自然な音符系列を推定する統計的手法を示す。歌声の発音時刻や F0 は楽譜に示されたビート時刻や音符の音高からの大きな逸脱を含むため、歌声 F0 軌跡の時間・周波数方向への離散化による音符推定の精度を向上するためには、楽譜の音楽的な自然さを表現する楽譜モデルが重要である。我々は調とリズムに依存する音符の音高を表現する楽譜モデルと楽譜 (音符系列) から時間・周波数方向に逸脱する歌声 F0 軌跡を表現する F0 モデルとを統合した階層隠れセミマルコフモデル (HHSMM: hierarchical hidden semi-Markov model) を提案する。楽譜モデルでは、確率的に生成された調に従って音符の音高が生成される。さらに、音符の開始位置はビートの 1 次元格子以上に定義されたマルコフ過程に従って生成される。F0 モデルでは、歌声の発音時刻の時間方向の逸脱、音符間における F0 の滑らかな遷移、F0 の周波数方向の逸脱が確率的に生成され、楽譜に付与される。提案法では、楽譜モデルと F0 モデルが音符推定に与える影響を考慮しながら、入力の歌声 F0 軌跡から尤もらしい音符系列を推定する。実験結果から調やリズムを考慮しない場合と比較して、提案法による音符系列の推定精度が向上することを示した。

キーワード 歌声, 自動採譜, 階層隠れセミマルコフモデル

1. はじめに

歌声は通常ポピュラー音楽のメロディラインを形成し、楽曲に関する多くの情報を提供するため、歌声解析は音楽情報検索をはじめとした様々な音楽アプリケーションにとって重要である。歌声 F0 推定 [1-5] や歌声分離 [6,7] といった歌声解析技術は盛んに研究されており、歌手同定 [8,9], カラオケ生成 [10], ハミング検索 [7], 能動的音楽鑑賞 [11] などに応用されている。さらに、歌声に含まれる情報をより活用するためには、歌声 F0 軌跡を離散的な記号のみを含む楽譜に変換することが有用である。

本研究では、楽譜に対して多くの逸脱を含む歌声 F0 軌跡からの音符系列推定に取り組む。楽譜における音符の音高や開始位置は離散的な値であるが、歌声 F0 軌跡は時間経過とともに滑らかに変化する連続的な信号である。例えば、歌声 F0 軌跡はビブラートによって振動したり、ポルタメントによってある音符から次の音符へと滑らかに変化したりする。したがって、歌声 F0 軌跡を単純に時間・周波数方向に離散化すると不自然なリズムや統計的に稀な半音階進行を含む音符系列がしばしば推定されてしまう。

この問題を解決するため、音符系列の生成過程を表現する楽譜モデルと歌声 F0 軌跡の生成過程を表現する F0 モデルの統合モデルに基づく調とリズムを考慮した統計的音符推定手法を

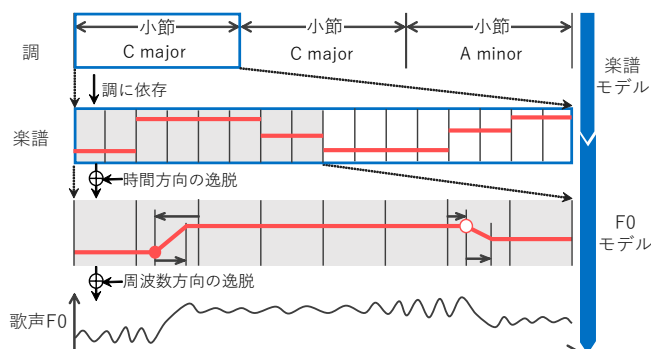


図 1: 楽譜モデルと F0 モデルの階層隠れセミマルコフモデルに基づく歌声 F0 軌跡の生成過程。

提案する (図 1)。楽譜モデルでは、各小節ごとに調がマルコフ過程により決定され、音符の音高は調と直前の音符の音高に依存して決まる。各調のもとでは、調の音階構成音に含まれる音名を持つ音高が出やすくなるよう制御される。ビートグリッド上に定義される各音符の開始位置は、直前の音符の開始位置に依存することで、リズム構造を形成する。F0 モデルでは、楽譜モデルによって生成された楽譜に対応する階段関数状の F0 軌跡に時間・周波数方向の逸脱が付与される。統合モデルは階層隠れセミマルコフモデル (HHSMM: hierarchical hidden semi-Markov model) として定式化される。提案法は歌声 F0

軌跡とビート時刻を入力として受け取り，マルコフ連鎖モンテカルロ法を用いることで，HSSMM の潜在変数として表現される調，音符，F0 の逸脱を同時に推定する．歌声 F0 軌跡の時間・周波数方向の離散化において調やリズムが自己組織化の制約として機能することが本手法の重要な特徴である．

2. 関連研究

本章では歌声解析に関する研究を紹介する．

2.1 音楽音響信号に対する歌声 F0 推定

音楽音響信号に対する歌声 F0 軌跡の推定は活発に研究されており [1–5]，これらの出力結果は提案法の入力として用いられる．最も基本的な方法の 1 つとして，各 F0 候補のそれぞれについて高調波成分の和を計算する Subharmonic Summation (SHS) [1] がある．池宮ら [2] は SHS に基づく歌声 F0 推定とロバスト主成分分析 (RPCA: robust principal component analysis) に基づく歌声分離の性能を，これら 2 つのタスクの相互依存性を利用することで改善した．Salamon ら [12] は特徴関数を計算することで歌声 F0 軌跡の候補を推定し，各軌跡の特徴から主旋律を形成しない軌跡を再帰的に消去する手法を提案した．Durrieu ら [3] は歌声と伴奏をそれぞれソース・フィルタモデルと非負値行列因子分解 (NMF: non-negative matrix factorization) に基づくモデルで表現することにより，主旋律の分離を行った．Mauch ら [5] は YIN [4] を確率的な手法に修正することで，システムが複数の F0 候補を出力し，その中から各フレームごとに 1 つの F0 を HMM を用いて選択するようにした．

2.2 歌声に対する音符推定

歌声に対する音符推定も盛んに研究が行われている [11, 13–19]．素朴な手法として，一定の区間ごとに歌声 F0 の多数決をとって音符の音高を決定する手法がある [11]．Paiva ら [13] は多重音検出，複数の F0 軌跡の構築，それら軌跡の分割，不要な音符の消去，主旋律を形成する音符の抽出の 5 つの処理を順番に行う手法を提案した．Raphael [14] は音符の個数を与えて，音高，リズム，テンポを推定する HMM に基づく手法を提案した．我々の提案法で用いられているリズムや歌声の発音時刻の逸脱に関するモデルは [14] で用いられたものと同様である．Laaksonen ら [15] は入力として与えたコードの境界に注目することで音響データを調と音符に対応する区間に分割し，スコア関数に基づいて各区間ごとに音符を推定する手法を提案した．Ryynänen ら [16] は 1 つの音符内における種々の歌声変動（例えば，ビブラートやポルタメント）を捉えるために階層 HMM に基づく手法を提案した．この手法のモデルでは，上層の HMM が音符の音高間の遷移を表し，下層の HMM が歌声変動の遷移を表す．Molina ら [17] は歌声 F0 軌跡における履歴現象に焦点を当てた．錦見ら [20] は時間・周波数方向の逸脱を考慮した歌声 F0 軌跡の生成過程を表現する HMM に基づく手法を提案した．Yang ら [18] は f_0 - Δf_0 平面の生成過程を表現する階層 HMM に基づく手法を提案した．Mauch ら [19] は音高抽出を行う Tony というソフトウェアツールを開発した．このツール内では，PYIN [5] を用いて歌声 F0 推定を

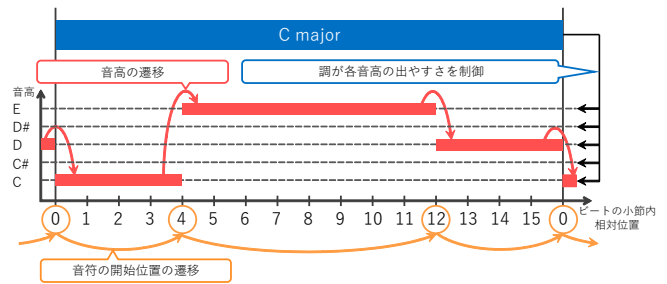


図 2: 楽譜モデル概要．

行い，Ryynänen's らの手法 [16] を基にした手法を用いて音符推定を行う．

3. 提案法

本章では歌声 F0 軌跡から音符系列を推定する提案法について説明する．提案法は，調に依存する音符系列から時間・周波数方向の逸脱を伴って歌声 F0 軌跡が確率的に生成される過程を HSSMM に基づいて表現する．提案モデルの上層は小節に割り当てられた調に従って音符系列が確率的に生成される過程を表現する HMM である．下層は時間方向の逸脱と周波数方向の逸脱がそれぞれ潜在変数と出力確率として表現される HSMM である．

3.1 問題設定

我々が取り組む問題を以下のように定める．

入力：歌声 F0 軌跡 $\mathbf{X} = \{x_t\}_{t=1}^T$ と 16 分音符単位のビート時刻

$\mathbf{Y} = \{(u_n, v_n)\}_{n=0}^N$ ，

出力：音符系列 $\mathbf{Z} = \{z_j = (p_j, l_j)\}_{j=0}^J$ ，

ここで， T は歌声 F0 軌跡のフレーム数， x_t は時刻 t における対数周波数， N は 16 分音符単位のビートの数である． $u_n \in \{1, \dots, T+1\}$ は n 番目のビート時刻であり，楽曲の最初と最後は $u_0 = 1$ と $u_N = T+1$ としてそれぞれ表される． $v_n \in \{0, \dots, 15\}$ は n 番目のビートが所属する小節内において，小節の先頭からそのビートまでの相対的な位置を表す． J は提案法によって推定される音符の個数であり， j 番目の音符 z_j は半音単位の音高 $p_j \in \{1, \dots, K\}$ と 16 分音符単位の音価 $l_j \in \{1, \dots, L\}$ の組として表現される．ここで， K は楽譜中に現れる音高の種類数であり， p_j は半音単位の音高に対応する対数周波数の集合 $\{\mu_1, \dots, \mu_K\}$ のうちの 1 つを指し示す．初期音符 z_0 は便宜上導入された実際の楽譜には現れない音符である．

3.2 楽譜の確率的モデル化

本章では音符の音高が調とリズムに依存して生成される過程を表現する HMM に基づく楽譜モデルについて説明する．

3.2.1 調遷移のモデル化

調系列は $\mathbf{S} = \{s_m\}_{m=0}^M$ で表現され， M は楽曲中の小節数， s_m は m 番目の小節の調を表す．便宜上，初期音符 z_0 が所属する初期小節を導入し，その小節に割り当てられる調を s_0 とする．転調に対応できるようにするため曲全体で調を 1 つに固定

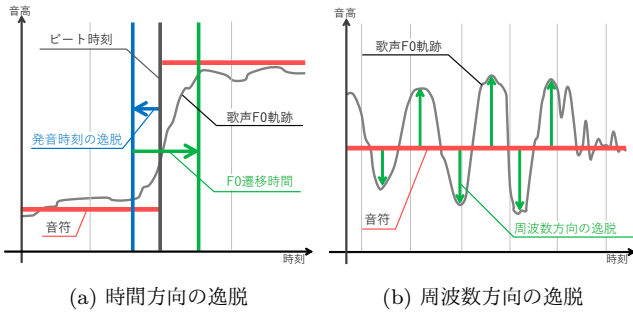


図 3: 歌声 F0 軌跡に含まれる逸脱。

せず、小節線で切り替わるようにする。各調 s_m は 24 通りの値 $\{C, C\#, \dots, B\} \times \{\text{major}, \text{minor}\}$ の中のいずれの値をとる。潜在変数 \mathbf{S} は以下のようにマルコフ連鎖をなす。

$$p(s_0 | \boldsymbol{\pi}) = \pi_{s_0} \quad (1)$$

$$p(s_m | s_{m-1}, \boldsymbol{\xi}_{s_{m-1}}) = \xi_{s_{m-1}s_m} \quad (2)$$

ここで、 $\boldsymbol{\pi} \in \mathbb{R}_{\geq 0}^{24}$ は初期確率、 $\boldsymbol{\xi}_s \in \mathbb{R}_{\geq 0}^{24}$ は遷移確率である。

3.2.2 音高遷移のモデル化

音高系列 \mathbf{P} は以下のように調系列 \mathbf{S} に依存したマルコフ連鎖によって生成される (図 2)。

$$p(p_0 | s_0, \boldsymbol{\phi}_{s_0}) = \phi_{s_0 p_0} \quad (3)$$

$$p(p_j | p_{j-1}, s_m, \boldsymbol{\psi}_{s_m p_{j-1}}) = \psi_{s_m p_{j-1} p_j} \quad (4)$$

ここで、 $\boldsymbol{\phi}_s \in \mathbb{R}_{\geq 0}^K$ は初期確率、 $\boldsymbol{\psi}_{sp} \in \mathbb{R}_{\geq 0}^K$ は遷移確率、 m は音符 z_j が属する小節のインデックスである。さらに、 $\phi_{s_0 p_0}$ と $\psi_{s_m p_{j-1} p_j}$ を以下のように定義する。

$$\phi_{s_0 p_0} = \frac{\hat{\phi}_{\hat{s}_0 \deg(p_0; s_0)}}{\sum_{p=1}^K \hat{\phi}_{\hat{s}_0 \deg(p; s_0)}} \quad (5)$$

$$\psi_{s_m p_{j-1} p_j} = \frac{\hat{\psi}_{\hat{s}_m \deg(p_{j-1}; s_m) \deg(p_j; s_m)}}{\sum_{p=1}^K \hat{\psi}_{\hat{s}_m \deg(p_{j-1}; s_m) \deg(p; s_m)}} \quad (6)$$

ここで、 $\hat{s} \in \{\text{major}, \text{minor}\}$ は調 s の旋法、 $\deg(p; s) \in \{0, \dots, 11\}$ は調 s における音高 p の度数 (調 s の主音に対する p のピッチクラスの音程) である。 $\hat{\phi}_*$ と $\hat{\psi}_*$ はそれぞれ旋法が与えられた下でのピッチクラスの初期確率と遷移確率である。

3.2.3 音符の開始位置遷移のモデル化

隣接する音符の開始位置間の遷移を考慮することで、音符系列 \mathbf{Z} が妥当なリズムを持つようにする。 j 番目の音符 z_j の開始位置を $r_{j-1} \in \{v_n\}_{n=1}^N$ とすると、音符の開始位置の遷移確率は以下のように与えられる。

$$p(r_j | r_{j-1}, \boldsymbol{\zeta}_{r_{j-1}}) = \zeta_{r_{j-1} r_j} \quad (7)$$

ここで、 r_{j-1} と r_j との間の距離は音符 z_j の音価 l_j となる。曲の最初と最後に関しては $r_0 = v_0$ と $r_J = v_N$ とする。

3.3 歌声 F0 軌跡の確率的モデル化

本章では歌声 F0 軌跡の生成過程を表現する HSM に基づく F0 モデルについて説明する。提案モデルでは、音符の音高

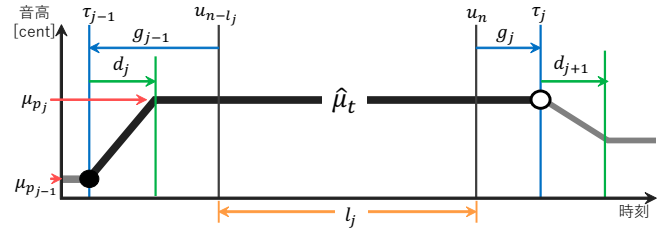


図 4: HSM の潜在変数と出力確率の位置パラメータの関係。黒い太線がコーシー分布の位置パラメータを表す。

と開始位置、時間方向の逸脱は潜在変数として表現され、周波数方向の逸脱は出力確率として表現する。

3.3.1 時間方向の逸脱のモデル化

歌声 F0 軌跡には以下のように 2 種類の時間方向の逸脱が含まれると仮定する (図 3a)。

発音時刻の逸脱: 歌声の発音時刻と音符の開始位置との間のずれ。

F0 の遷移時間: ある音符の音高から次の音符の音高まで、歌声が遷移し切るのに要する時間。

音符系列 \mathbf{Z} に付与される発音時刻の逸脱 $\mathbf{G} = \{g_j\}_{j=0}^J$ は離散潜在変数として表現される。音符の開始位置モデルと同様に音符 z_j の発音時刻の逸脱を g_{j-1} とする。各 g_j は $-G$ から G までの整数値を取り、以下のようにそれぞれ独立に生成されるとする。

$$p(g_j | \boldsymbol{\rho}) = \rho_{g_j} \quad (8)$$

ここで、 $\boldsymbol{\rho} \in \mathbb{R}_{\geq 0}^{2G+1}$ は発音時刻逸脱の確率の集合である。また、最初の音符の開始時刻と最後の音符の終了時刻には逸脱が無い、すなわち $g_0 = g_J = 0$ であるとする。

音符系列 \mathbf{Z} に付与される F0 の遷移時間 $\mathbf{D} = \{d_j\}_{j=1}^J$ は離散潜在変数として表現され、各 d_j は 1 から D までの整数値をとる。音符 z_{j-1} と z_j の間における歌声 F0 軌跡の連続的な遷移は、幅が d_j フレームの斜め線によって表現される。各 d_j は以下のように独立に生成される。

$$p(d_j | \boldsymbol{\eta}) = \eta_{d_j} \quad (9)$$

ここで、 $\boldsymbol{\eta} \in \mathbb{R}_{\geq 0}^D$ は F0 の遷移時間の確率の集合である。

3.3.2 周波数方向の逸脱のモデル化

歌声 F0 軌跡 $\mathbf{X} = \{x_t\}_{t=1}^T$ は時間方向の逸脱が既に付与された音符系列に対して周波数方向の逸脱がさらに付与されて生成される (図 3b)。 x_t は各フレームごとに独立に生成されるとし、 j 番目の音符 z_j に関する出力確率は以下の通りである。

$$\begin{aligned} & p(x_{\tau_{j-1}:\tau_j-1} | p_{j-1}, p_j, l_j, g_{j-1}, g_j, d_j, \hat{\mu}_t, \lambda) \\ &= \prod_{t=\tau_{j-1}}^{\tau_j-1} \{\delta_{x_t, \text{voiced}} \text{Cauchy}(x_t | \hat{\mu}_t, \lambda) + \delta_{x_t, \text{unvoiced}}\} \\ &= e^{p_{j-1} p_j l_j g_{j-1} g_j d_j} \end{aligned} \quad (10)$$

ここで、 $x_{\tau:\tau-1}$ は $x_\tau, \dots, x_{\tau-1}$ を表し、 λ は周波数方向の逸脱の大きさを表現する尺度パラメータ、 δ はクロネッカーのデ

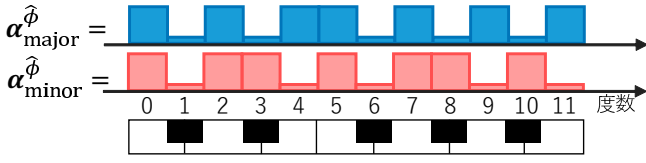


図 5: ハイパーパラメータ $\alpha_s^{\hat{\phi}}$ の設定.

ルタ, $\hat{\mu}_t$ (図 4) は以下のように定義される位置パラメータである.

$$\hat{\mu}_t = \begin{cases} \frac{\mu_{p_j} - \mu_{p_{j-1}}}{d_j} (t - \tau_{j-1}) + \mu_{p_{j-1}} & (\tau_{j-1} \leq t < \tau_j + d_j) \\ \mu_{p_j} & (\tau_{j-1} + d_j \leq t < \tau_j) \end{cases} \quad (11)$$

音符 z_{j+1} の開始位置が n 番目のビートに位置しているとき, $\tau_j = u_n + g_j$, $\tau_{j-1} = u_{n-l_j} + g_{j-1}$ である.

3.4 事前分布の導入

離散モデルパラメータ $\pi, \xi, \hat{\phi}, \hat{\psi}, \zeta, \rho, \eta$ に対して以下のようにディリクレ共役事前分布を置く.

$$\begin{aligned} \pi &\sim \text{Dirichlet}(\mathbf{a}^\pi) & \xi_s &\sim \text{Dirichlet}(\mathbf{a}_s^\xi) \\ \hat{\phi}_s &\sim \text{Dirichlet}(\mathbf{a}_s^{\hat{\phi}}) & \hat{\psi}_{\text{sdeg}(p;s)} &\sim \text{Dirichlet}(\mathbf{a}_{\text{sdeg}(p;s)}^{\hat{\psi}}) \\ \zeta_r &\sim \text{Dirichlet}(\mathbf{a}_r^\zeta) \\ \rho &\sim \text{Dirichlet}(\mathbf{a}^\rho) & \eta &\sim \text{Dirichlet}(\mathbf{a}^\eta) \end{aligned} \quad (12)$$

ここで, $\mathbf{a}^\pi \in \mathbb{R}_+^{26}$, $\mathbf{a}_s^\xi \in \mathbb{R}_+^{26}$, $\mathbf{a}_s^{\hat{\phi}} \in \mathbb{R}_+^{12}$, $\mathbf{a}_{\text{sdeg}(p;s)}^{\hat{\psi}} \in \mathbb{R}_+^{12}$, $\mathbf{a}_r^\zeta \in \mathbb{R}_+^{16}$, $\mathbf{a}^\rho \in \mathbb{R}_+^{2G+1}$, $\mathbf{a}^\eta \in \mathbb{R}_+^D$ はハイパーパラメータである. ある調のもとでの各ピッチクラスの出やすさは, それらピッチクラスの初期確率と遷移確率に関する事前分布を用いて制御される. 図 5 に示すように, ハイパーパラメータ $\mathbf{a}_s^{\hat{\phi}}$ と $\mathbf{a}_{\text{sdeg}(p;s)}^{\hat{\psi}}$ はそれぞれダイアトニックスケールを表現するように設定される. コーシー分布は共役事前分布を持たないので, 尺度パラメータ λ に対して以下のようにガンマ事前分布を置く.

$$\lambda \sim \text{Gamma}(a_0^\lambda, a_1^\lambda) \quad (13)$$

ここで, a_0^λ と a_1^λ はハイパーパラメータである.

3.5 ベイズ推定

我々の目的は歌声 F0 軌跡 \mathbf{X} が与えられた下で事後分布 $p(\mathbf{S}, \mathbf{Q}, \Theta | \mathbf{X})$ を計算することである. ここで, $\mathbf{Q} = \{\mathbf{P}, \mathbf{L}, \mathbf{G}, \mathbf{D}\}$ (潜在変数), $\Theta = \{\pi, \xi, \hat{\phi}, \hat{\psi}, \zeta, \rho, \eta, \lambda\}$ (モデルパラメータ) である. この事後分布は解析的に計算することが困難であるため, 我々はマルコフ連鎖モンテカルロ法 (MCMC: Markov chain Monte Carlo) を用いて $\mathbf{S}, \mathbf{Q}, \Theta$ の値をサンプルする. 潜在変数 \mathbf{S} と \mathbf{Q} のサンプルにはフォワードフィルタリング・バックワードサンプリングアルゴリズムを用いる. モデルパラメータ Θ のうち λ 以外, すなわち共役事前分布を持つパラメータのサンプルにはギブスサンプリングアルゴリズムを用いる. パラメータ λ には共役事前分布が無いので, メトロポリス・ヘイスティングス (MH: Metropolis-Hastings) アルゴリズムを用いる. \mathbf{S} と \mathbf{Q} は音符系列 \mathbf{Z} を共有し, 相互に依存しているため, 各変数は以下の手順で更新される.

- (1) 多数決法により音符系列 \mathbf{Z} を初期化する.

- (2) \mathbf{Z} に基づき調系列 \mathbf{S} を更新する.
- (3) \mathbf{S} に基づき \mathbf{Q} を更新する.
- (4) モデルパラメータ Θ を更新する.
- (5) 2 に戻る.

3.5.1 潜在変数 \mathbf{S} の推論

音符系列 \mathbf{Z} が与えられた下で各 s_m は以下に示す確率に従いサンプルされる.

$$\beta_{s_m}^{\mathbf{S}} = p(s_m | s_{m+1:M}, \mathbf{Z}) \quad (14)$$

ここで, $s_{m+1:M}$ は s_{m+1}, \dots, s_M を表す. 式 (14) の計算と調 \mathbf{S} のサンプルにはフォワードフィルタリング・バックワードサンプリング法を用いる.

フォワードフィルタリングでは確率 $\alpha_{s_m}^{\mathbf{S}}$ が以下のように再帰的に計算される.

$$\alpha_{s_0}^{\mathbf{S}} = p(p_0, s_0) = p(p_0 | s_0) p(s_0) = \phi_{s_0 p_0} \pi_{s_0} \quad (15)$$

$$\begin{aligned} \alpha_{s_m}^{\mathbf{S}} &= p(p_0 : j_{m+1}-1, s_m) \\ &= \prod_{j=j_m}^{j_{m+1}-1} \psi_{s_m p_{j-1} p_j} \sum_{s_{m-1}} \xi_{s_{m-1} s_m} \alpha_{s_{m-1}}^{\mathbf{S}} \end{aligned} \quad (16)$$

ここで, j_m は m 番目の小節に属する最初の音符のインデックスであり, 既知の音価 \mathbf{L} から計算できる.

バックワードサンプリングではフォワードフィルタリングで計算された値を用いて式 (14) を計算し, 以下のように調が再帰的にサンプルされる.

$$\beta_{s_M}^{\mathbf{S}} = p(s_M | \mathbf{Z}) \propto \alpha_{s_M}^{\mathbf{S}} \quad (17)$$

$$\beta_{s_m}^{\mathbf{S}} = p(s_m | s_{m+1:M}, \mathbf{Z}) \propto \alpha_{s_m}^{\mathbf{S}} \xi_{s_m s_{m+1}} \quad (18)$$

3.5.2 潜在変数 \mathbf{Q} の推論

潜在変数 \mathbf{Q} は \mathbf{S} と同様の方法で推論される. フォワードフィルタリングでは確率 $\alpha_{p_n l_n g_n d_n}^{\mathbf{Q}}$ が以下のように再帰的に計算される.

$$\alpha_{p_0 l_0 g_0 d_0}^{\mathbf{Q}} = p(p_0 | \mathbf{S}) = \phi_{y_0 p_0} \quad (19)$$

$$\alpha_{p_n l_n g_n d_n}^{\mathbf{Q}} = p(x_{1:\tau_n-1}, p_n, l_n, g_n, d_n | \mathbf{S}) = \begin{cases} 0 & (l_n > n) \\ \rho_{g_n} \eta_{d_n} \zeta_{r_0 r_n} \\ \cdot \sum_{p_0} \psi_{s_1 p_0 p_n} e_{p_0 p_n l_n g_n d_n} \alpha_{p_0 l_0 g_0 d_0}^{\mathbf{Q}} & (l_n = n) \\ \sum_{p_{n'}: g_{n'}} \sum_{l_{n'}} \sum_{d_{n'}} \rho_{g_n} \eta_{d_n} \zeta_{r_{n'} r_n} \psi_{s_{m(n')} p_{n'} p_n} \\ \cdot e_{p_{n'} p_n l_{n'} g_{n'} d_{n'}} \alpha_{p_{n'} l_{n'} g_{n'} d_{n'}}^{\mathbf{Q}} & (l_n < n) \end{cases} \quad (20)$$

ここで, $\tau_n = u_n + g_n$, $n' = n - l_n$ であり, $m(n')$ は n' 番目のビートが属する小節のインデックスである. p_n, l_n, g_n, d_n は終了位置が n 番目のビート u_n に位置する音符に対応するフォワードメッセージの変数である. これらの変数は j を添字とする変数 p_j, l_j, g_j, d_j とは異なる. 音符 $z_n = (p_n, l_n)$ の開始位置と終了位置はそれぞれ $(n-l_n)$ 番目のビートと n 番目のビートに位置し, 式 (20) の再帰計算に現れる確率 $p(l_n)$ は

$p(r_n|r_{n-l_n})$ に置き換えられる。

バックワードサンプリングではフォワードフィルタリングで計算された値を用いて潜在変数の事後分布を計算し、以下のよう
に音符と時間方向の逸脱が再帰的にサンプルされる。

$$\begin{aligned} \beta_{p_N l_N g_N d_N} &= p(p_N, l_N, g_N, d_N | \mathbf{X}, \mathbf{S}) \propto \alpha_{p_N l_N g_N d_N}^Q \\ \beta_{p_{n'} l_{n'} g_{n'} d_{n'}} &= p(p_{n'}, l_{n'}, g_{n'}, d_{n'} | p_{n:N}, l_{n:N}, g_{n:N}, d_{n:N}, \mathbf{X}) \\ &\propto \begin{cases} 0 & (l_n > n) \\ e_{p_{n'} p_n l_{n'} g_{n'} d_{n'}} \psi_{s_{m(n')} p_{n'} p_n} \\ \cdot \zeta_{r_{n'} r_n} \rho_{g_n} \eta_{d_n} \alpha_{p_{n'} l_{n'} g_{n'} d_{n'}}^Q & (l_n \leq n) \end{cases} \end{aligned} \quad (21)$$

3.5.3 モデルパラメータ Θ の学習

共役事前分布を持つモデルパラメータの事後分布はバック
ワードサンプリングで得られたサンプル \mathbf{S}, \mathbf{Q} を用いて計算さ
れる。そして、これらのパラメータは計算された事後分布に従
い以下のようにサンプルされる。

$$\boldsymbol{\pi} \sim \text{Dirichlet}(\mathbf{a}^\pi + \mathbf{b}^\pi) \quad \boldsymbol{\xi}_s \sim \text{Dirichlet}(\mathbf{a}_s^\xi + \mathbf{b}_s^\xi) \quad (22)$$

$$\hat{\boldsymbol{\phi}}_s \sim \text{Dirichlet}(\mathbf{a}_s^{\hat{\phi}} + \mathbf{b}_s^{\hat{\phi}}) \quad (23)$$

$$\hat{\boldsymbol{\psi}}_{\text{sdeg}(p;s)} \sim \text{Dirichlet}(\mathbf{a}_{\text{sdeg}(p;s)}^{\hat{\psi}} + \mathbf{b}_{\text{sdeg}(p;s)}^{\hat{\psi}}) \quad (24)$$

$$\boldsymbol{\zeta}_r \sim \text{Dirichlet}(\mathbf{a}_r^\zeta + \mathbf{b}_r^\zeta) \quad (25)$$

$$\boldsymbol{\rho} \sim \text{Dirichlet}(\mathbf{a}^\rho + \mathbf{b}^\rho) \quad \boldsymbol{\eta} \sim \text{Dirichlet}(\mathbf{a}^\eta + \mathbf{b}^\eta) \quad (26)$$

ここで、 $\mathbf{b}^\pi \in \mathbb{R}_{\geq 0}^{26}$ は第 s_0 成分だけが 1 の単位ベクトルである。
 $\mathbf{b}_s^\xi \in \mathbb{R}_{\geq 0}^{26}$ はサンプル \mathbf{Y} における調 s から s' への遷移の回数を
第 s' 成分が表すベクトルである。 $\mathbf{b}^\rho \in \mathbb{R}_{\geq 0}^{2G+1}$ はサンプル \mathbf{Q} に
おける歌声の発音時刻の逸脱が g である回数を第 g 成分が表す
ベクトルであり、 $\mathbf{b}^\eta \in \mathbb{R}_{\geq 0}^D$ はサンプル \mathbf{Q} における F0 の遷移時
間 d である回数を第 d 成分が表すベクトルである。 $\mathbf{b}_r^\zeta \in \mathbb{R}_{\geq 0}^{16}$
は、バックワードサンプリングで得られた音価 \mathbf{L} から計算され
る $\mathbf{R} = \{r_j\}_{j=0}^J$ において、音符の開始位置 r から r' への遷移の
回数を第 r' 成分が表すベクトルである。ベクトル $\mathbf{b}_s^{\hat{\phi}} \in \mathbb{R}_{> 0}^{12}$ は、
初期小節の調と初期音符の音高がそれぞれ $s_0 = s, p_0 = p$ で
ある時、要素 $b_{\text{sdeg}(p;s)}^{\hat{\phi}}$ の値が 1 でそれ以外の要素の値が 0 で
あるようなベクトルである。ベクトル $\mathbf{b}_{\text{sdeg}(p;s)}^{\hat{\psi}} \in \mathbb{R}_{\geq 0}^{12}$ は、サン
プルされた潜在変数において調 s の下で音高 p から p' への遷
移の回数を要素 $b_{\text{sdeg}(p;s)\text{deg}(p';s)}^{\hat{\psi}}$ が表すベクトルである。

尺度パラメータ λ については MH アルゴリズムを適用する
ため、以下のように提案分布を定める。

$$q(\lambda^* | \lambda) = \text{Gamma}(\gamma \lambda, \gamma) \quad (27)$$

ここで、 λ^* は次の λ の値の候補を表す変数、 λ は現在の尺度パ
ラメータの値、 γ はハイパーパラメータである。 λ の値の更新
は以下の確率に従い行われる。

$$A(\lambda^*, \lambda) = \min \left\{ \frac{\mathcal{L}(\lambda^*) q(\lambda | \lambda^*)}{\mathcal{L}(\lambda) q(\lambda^* | \lambda)} \right\} \quad (28)$$

ここで、 $L(\lambda)$ は以下のように与えられる。

$$\mathcal{L}(\lambda) = \text{Gamma}(\lambda | a_0^\lambda, a_1^\lambda) \prod_{j=1}^J e_{p_{j-1} p_j l_j g_j d_j} \quad (29)$$

$\{p_j, l_j, g_j, d_j\}_{j=0}^J$ はバックワードサンプリングでサンプルされ
た値である。一様分布 $\mathcal{U}(0, 1)$ からサンプルされた乱数よりも
採択率 $A(\lambda^*, \lambda)$ の値が大きい場合に λ の値は λ^* に更新される。

3.6 ビタビ復号

潜在変数系列 \mathbf{S}, \mathbf{Q} は学習時に同時分布 $p(\mathbf{X}, \mathbf{Q}, \mathbf{S}, \Theta | \Phi)$ を
最大化したモデルパラメータの値を用いるビタビアルゴリズム
によって推定される。潜在変数の推論と同様に、多数決法によ
り音符系列 \mathbf{Z} を初期化したのち、 \mathbf{Z} に基づき \mathbf{S} を推定し、推
定した \mathbf{S} に基づき \mathbf{Q} を推定する。

\mathbf{S} に関するビタビ復号では、以下のように ω_s^S が再帰的に計
算される。

$$\omega_{s_0}^S = \ln \phi_{s_0 k_0} + \ln \pi_{s_0} \quad (30)$$

$$\omega_{s_m}^S = \sum_{j=j_m}^{j_{m+1}-1} \ln \psi_{s_m p_{j-1} p_j} + \max_{s_{m-1}} \{ \ln \xi_{s_{m-1} s_m} + \omega_{s_{m-1}}^S \} \quad (31)$$

$\omega_{s_m}^S$ の再帰計算では、 $\omega_{s_m}^S$ の値を最大化する 1 つ前の状態 s_{m-1}
が $c_{s_m}^S$ として記録され、調系列 \mathbf{S} は以下のように再帰的に推
定される。

$$s_M = \arg \max_{s_M} \alpha_{s_M}^S \quad s_{m-1} = c_{s_m}^S \quad (32)$$

\mathbf{Q} に関するビタビ復号では、 ω_{plgd}^Q の値が以下のように再帰
的に計算される。

$$\omega_{p_0 l_0 g_0 d_0}^Q = w^\phi \ln \phi_{s_0 p_0} \quad (33)$$

$$\omega_{p_n l_n g_n d_n}^Q = \begin{cases} -\text{inf} & (l_n > n) \\ w^\rho \ln \rho_{g_n} + w^\eta \ln \eta_{d_n} + w^\zeta \ln \zeta_{r_n r_n} \\ + \max_{p_0} \left\{ w^\psi \ln \psi_{s_1 p_0 p_n} \right. \\ \left. + w^e \ln e_{p_0 p_n l_n g_n d_n} + \omega_{p_0 l_0 g_0 d_0}^Q \right\} & (l_n = n) \\ w^\rho \ln \rho_{g_n} + w^\eta \ln \eta_{d_n} + w^\zeta \ln \zeta_{r_n r_n} \\ + \max_{(p_{n'}, l_{n'}, g_{n'}, d_{n'})} \left\{ w^\psi \ln \psi_{s_{m(n')} p_{n'} p_n} \right. \\ \left. + w^e \ln e_{p_{n'} p_n l_{n'} g_{n'} d_{n'}} + \omega_{p_{n'} l_{n'} g_{n'} d_{n'}}^Q \right\} & (l_n < n) \end{cases} \quad (34)$$

ここで、 $w^\phi, w^\psi, w^\rho, w^\eta, w^\zeta, w^e$ は各確率間のバランスを制御
する重みパラメータである。 ω_{plgd}^Q の再帰計算では、 $\omega_{p_n l_n g_n d_n}^Q$
の値を最大化する 1 つ前の状態 $p_{n'}, l_{n'}, g_{n'}, d_{n'}$ が $c_{p_n l_n g_n d_n}^Q$
として記録され、 \mathbf{Q} は以下のように再帰的に推定される。

$$(p_N, l_N, g_N, d_N) = \arg \max_{p_N, l_N, g_N, d_N} \alpha_{p_N l_N g_N d_N}^Q \quad (35)$$

$$(p_{n'}, l_{n'}, g_{n'}, d_{n'}) = c_{p_n l_n g_n d_n}^Q \quad (36)$$

4. 評価実験

歌声 F0 軌跡からの音符推定について、提案法の精度を評価

表 1: ビート単位と音符単位の一致率 [%] および標準誤差

モデル	入力 F0	ビート単位	音符単位
提案法	正解	72.4 ± 1.7	28.1 ± 2.1
	推定	68.7 ± 1.3	30.7 ± 1.8
リズムのみ考慮	正解	71.5 ± 1.6	26.3 ± 2.1
	推定	67.7 ± 1.3	29.1 ± 1.8
調のみ考慮	正解	67.8 ± 1.6	10.6 ± 1.2
	推定	65.6 ± 1.2	13.8 ± 1.1
調・リズムを考慮しない	正解	67.2 ± 1.5	9.8 ± 1.2
	推定	64.6 ± 1.2	12.9 ± 1.1
多数決法	正解	54.1 ± 1.5	20.1 ± 1.4
	推定	61.0 ± 1.4	22.0 ± 1.5
HMM [20]	推定	68.0 ± 1.2	14.8 ± 1.3

するために比較実験を行った。

4.1 実験条件

RWC 研究用音楽データベース [21] のポピュラー音楽 100 曲のうち、提案法が扱えない 32 分音符、3 連符、ハモリパートなどを含む曲を除いた 63 曲を用いた。入力の歌声 F0 軌跡 X はアノテーションデータ [22] と [2] で提案されている手法によって推定されたものを用いた。アノテーションデータには無声区間が含まれるが推定データには含まれない。ビート時刻とビートの小節内における相対位置 Y はアノテーションデータから得た。

ベイズ推定とビタビ復号は各曲独立に行われる。音符の開始位置の遷移確率はロック音楽のコーパス [23] から事前に学習した。ハイパーパラメータは $\alpha^\pi=1, \alpha_s^\xi=1, \alpha_s^\zeta=1, \alpha^p = \alpha^\eta = \alpha_0^\lambda = \alpha_1^\lambda = \gamma = 1$, とした。ここで、 $\mathbb{1}$ と $\mathbf{1}$ はそれぞれすべての要素が 1 の行列とベクトルである。 α_s^ϕ と $\alpha_{s^{\text{deg}(p;s)}}^\psi$ は、調 s のダイアトニックスケール内のピッチクラスに対応する要素を 10、それ以外の要素を 1 に設定した。ビタビアルゴリズムの重みパラメータは経験的に $w^\phi = w^\psi = 29.4, w^p = 2.4, w^\eta = 2.9, w^\zeta = 48.5, w^e = 3.8$ とした。音楽的に一貫性のある音符系列を得るために、F0 モデルに関する重みパラメータよりも楽譜モデルに関する重みパラメータの値を大きくした。

比較ために、多数決法と HMM に基づく従来法 [20] についても実験した。従来法 [20] については、無声区間を含む歌声 F0 軌跡を入力として扱えないため、無声区間を含まない推定データについてのみ実験した。言語モデルの有効性を評価するために、提案法に関しても、1) 調とリズムのどちらも考慮しない手法、2) 調のみを考慮した手法、3) リズムのみを考慮した手法、4) 調とリズムの両方を考慮した手法、の 4 通りの実験を行った。提案法における学習を高速化するため、音符の音高の探索範囲を多数決法によって推定された音高の周囲に限定した。

各手法の性能を評価するため、正解の音符系列と推定された音符系列を比較して、ビート単位の一一致率と音符単位の一一致率を計算した。ビート単位の一一致率は正解の楽譜内の音符が存在するビート区間の個数に対して、正しく音高が推定されたビート区間の個数の割合とする。音符単位の一一致率は正解楽譜内の

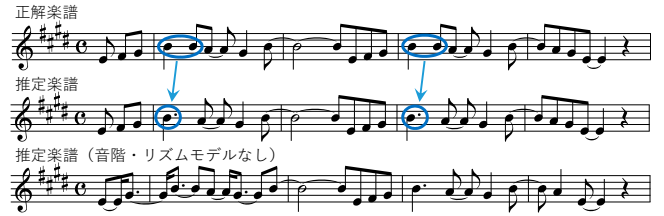


図 6: 提案法と調とリズムのどちらも考慮しない手法を用いて正解の歌声 F0 軌跡から推定された楽譜の例。

音符の個数に対して、音高、開始位置、終了位置の 3 つ全てが正しく推定された音符の個数の割合とする。正解楽譜の中の隣接する音符が、同じ音高であるかタイで結ばれている場合、それらの音符をまとめて 1 つの音符とみなした。従来法 [20] では 16 分音符ごとに音高を出力するので、連続する同じ音高の並びを 1 つの音符とみなした。

4.2 実験結果

実験結果を表 1 に示した。提案法は両評価尺度において、多数決法や従来法よりも音符推定精度が上回っていた。提案法に関する 4 通りの実験から得られたビート単位の一一致率を比較すると、楽譜モデルによって音符推定の性能が向上することを確かめた。特に、調の遷移確率 (調の制約) よりも音符の開始時刻の遷移確率 (リズムの制約) の方がより有効であることが分かった。ビート単位の一一致率では、提案法 (68.7%) と従来法 (68.0%) で大きな差は見られなかったが、音符単位の一一致率では、提案法が (30.7%) が従来法 (14.8%) を大きく上回った。

推定された楽譜の例を図 6 に示す。一部の音符が結合されたことを除いて、提案法により推定された楽譜がほぼ正確であることがわかる。隣接する同じ音高の音符を正しく推定するには歌声 F0 軌跡だけでは限界があり、元の歌声や音楽音響信号を参照する必要がある。一方で、楽譜モデルを考慮せずに推定された楽譜には多くの推定誤りが含まれていた。この結果からも、音符推定において楽譜モデルを音楽的制約として用いることの有効性が示せた。

5. おわりに

本稿では歌声 F0 軌跡から音符系列を推定する統計的手法を示した。提案法は調から楽譜が生成される過程を表す楽譜モデルと楽譜から時間・周波数方向の逸脱を伴って歌声 F0 軌跡が生成される過程を表す F0 モデルを統合した階層隠れセミマルコフモデル (HHSMM) に基づいており、音楽的に一貫性のある音符系列を出力できることを確かめた。

本研究の今後の方向として最も興味深いのは、音楽音響信号に対する歌声 F0 推定において歌声 F0 軌跡の音楽的に有意な事前分布として提案したモデルを用いることである。本稿で提案した楽譜から F0 軌跡を生成するモデルを「言語」モデルとし、歌声 F0 軌跡からスペクトログラムを出力する音響モデルと階層ベイズの枠組みで統合する予定である。これにより音楽音響信号から歌声 F0 軌跡と楽譜を同時に学習することが可能になる。また、提案法では事前に推定した歌声 F0 軌跡とビート時

刻を入力として与えているが，入力の推定精度が音符推定精度に影響を与える問題を克服するためにも歌声 F0 軌跡とビート時刻の同時推定法について検討すべきである。

謝辞：本研究の一部は，JSPS 科研費 26700020, 16H01744, 16J05486 および JST ACCEL No. JPMJAC1602 の支援を受けた。

文 献

- [1] D.J. Hermes, “Measurement of pitch by subharmonic summation,” *The journal of the acoustical society of America*, vol.83, no.1, pp.257–264, 1988.
- [2] Y. Ikemiya, K. Yoshii, and K. Itoyama, “Singing voice analysis and editing based on mutually dependent F0 estimation and source separation,” 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2015), pp.574–578, 2015.
- [3] J.-L. Durrieu, G. Richard, B. David, and C. Févotte, “Source/filter model for unsupervised main melody extraction from polyphonic audio signals,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol.18, no.3, pp.564–575, 2010.
- [4] A. deCheveigné and H. Kawahara, “YIN, a fundamental frequency estimator for speech and music,” *The Journal of the Acoustical Society of America*, vol.111, no.4, pp.1917–1930, 2002.
- [5] M. Mauch and S. Dixon, “pYIN: A fundamental frequency estimator using probabilistic threshold distributions,” 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2014), pp.659–663, 2014.
- [6] Y. Li and D. Wang, “Separation of singing voice from music accompaniment for monaural recordings,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol.15, no.4, pp.1475–1487, 2007.
- [7] P.-S. Huang, S.D. Chen, P. Smaragdis, and M. Hasegawa-Johnson, “Singing-voice separation from monaural recordings using robust principal component analysis,” 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2012), pp.57–60, 2012.
- [8] Y.E. Kim and B. Whitman, “Singer identification in popular music recordings using voice coding features,” 3rd International Conference on Music Information Retrieval (ISMIR 2002), vol.13, p.17, 2002.
- [9] W.-H. Tsai and H.-M. Wang, “Automatic singer recognition of popular music recordings via estimation and modeling of solo vocal signals,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol.14, no.1, pp.330–341, 2006.
- [10] M. Ryyänänen, T. Virtanen, J. Paulus, and A. Klapuri, “Accompaniment separation and karaoke application based on automatic melody transcription,” 2008 IEEE International Conference on Multimedia and Expo, pp.1417–1420, 2008.
- [11] M. Goto, K. Yoshii, H. Fujihara, M. Mauch, and T. Nakano, “Songle: A web service for active music listening improved by user contributions,” *Proc. of the 12th International Society for Music Information Retrieval Conference (ISMIR 2011)*, pp.311–316, 2011.
- [12] J. Salamon and E. Gómez, “Melody extraction from polyphonic music signals using pitch contour characteristics,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol.20, no.6, pp.1759–1770, 2012.
- [13] R.P. Paiva, T. Mendes, and A. Cardoso, “On the detection of melody notes in polyphonic audio,” 6th International Conference on Music Information Retrieval (ISMIR 2005), pp.175–182, 2005.
- [14] C. Raphael, “A graphical model for recognizing sung melodies,” 6th International Conference on Music Information Retrieval (ISMIR 2005), pp.658–663, 2005.
- [15] A. Laaksonen, “Automatic melody transcription based on chord transcription,” *Proc. of the 15th International Society for Music Information Retrieval (ISMIR 2014)*, pp.119–124, 2014.
- [16] M.P. Ryyänänen and A.P. Klapuri, “Automatic transcription of melody, bass line, and chords in polyphonic music,” *Computer Music Journal*, vol.32, no.3, pp.72–86, 2008.
- [17] E. Molina, L.J. Tardón, A.M. Barbancho, and I. Barbancho, “Sipth: Singing transcription based on hysteresis defined on the pitch-time curve,” *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)*, vol.23, no.2, pp.252–263, 2015.
- [18] L. Yang, A. Maezawa, J.B.L. Smith, and E. Chew, “Probabilistic transcription of sung melody using a pitch dynamic model,” 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2017), pp.301–305, 2017.
- [19] M. Mauch, C. Cannam, R. Bittner, G. Fazekas, J. Salamon, J. Dai, J. Bello, and S. Dixon, “Computer-aided melody note transcription using the Tony software: Accuracy and efficiency,” *Proc. of the 1st International Conference on Technologies for Music Notation and Representation (TENOR 2015)*, pp.23–30, 2015.
- [20] R. Nishikimi, E. Nakamura, K. Itoyama, and K. Yoshii, “Musical note estimation for f0 trajectories of singing voices based on a bayesian semi-beat-synchronous hmm,” *Proc. of the 17th International Society for Music Information Retrieval Conference (ISMIR 2016)*, pp.461–467, 2016.
- [21] M. Goto, H. Hashiguchi, T. Nishimura, and R. Oka, “RWC music database: Popular, classical and jazz music databases,” *The 3rd International Conference on Music Information Retrieval (ISMIR 2002)*, pp.287–288, 2002.
- [22] M. Goto, “Aist annotation for the RWC music database,” *The 7th International Conference on Music Information Retrieval (ISMIR 2006)*, pp.359–360, 2006.
- [23] T. De Clercq and D. Temperley, “A corpus analysis of rock harmony,” *Popular Music*, vol.30, no.01, pp.47–70, 2011.