

SCALE- AND RHYTHM-AWARE MUSICAL NOTE ESTIMATION FOR VOCAL F0 TRAJECTORIES BASED ON A SEMI-TATUM-SYNCHRONOUS HIERARCHICAL HIDDEN SEMI-MARKOV MODEL

Ryo Nishikimi¹ Eita Nakamura¹ Masataka Goto² Katsutoshi Itoyama¹ Kazuyoshi Yoshii^{1,3}

¹Graduate School of Informatics, Kyoto University, Japan ³RIKEN AIP, Japan

²National Institute of Advanced Industrial Science and Technology (AIST), Japan

{nishikimi, enakamura, itoyama, yoshii}@sap.ist.i.kyoto-u.ac.jp, m.goto@aist.go.jp

ABSTRACT

This paper presents a statistical method that estimates a sequence of musical notes from a vocal F0 trajectory. Since the onset times and F0s of sung notes are considerably deviated from the discrete tatums and pitches indicated in a musical score, a score model is crucial for improving time-frequency quantization of the F0s. We thus propose a hierarchical hidden semi-Markov model (HHSMM) that combines a score model representing the rhythms and pitches of musical notes with musical scales with an F0 model representing time-frequency deviations from a note sequence specified by a score. In the score model, musical scales are generated stochastically. Note pitches are then generated according to the scales and note onsets are generated according to a Markov process defined on the tatum grid. In the F0 model, onset deviations, smooth note-to-note F0 transitions, and F0 deviations are generated stochastically and added to the note sequence. Given an F0 trajectory, our method estimates the most likely sequence of musical notes while giving more importance on the score model than the F0 model. Experimental results showed that the proposed method outperformed an HMM-based method having no models of scales and rhythms.

1. INTRODUCTION

Singing voice analysis is important for music information retrieval because a singing voice usually forms a large part of the melody line of popular music, and provides much information about music. Singing voice analysis techniques such as vocal F0 estimation [1, 3, 7, 9, 14] and singing voice separation [8, 12] have actively been studied and applied to singer identification [10, 22], karaoke generation [19], query-by-humming [8], and active music listening [6]. To leverage musical information conveyed by singing voices, it is helpful to convert a vocal F0 trajectory to a musical score containing only discrete symbols.

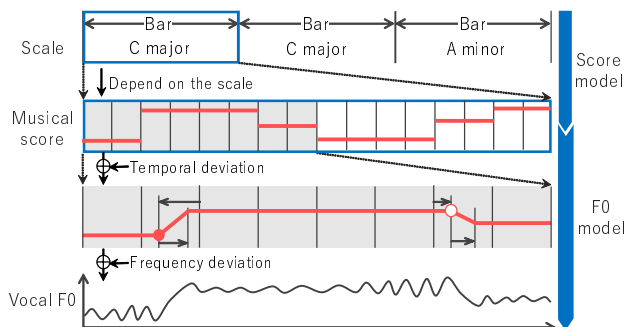


Figure 1: The generative process of a vocal F0 trajectory based on a hierarchical hidden semi-Markov model involving a score model and an F0 model.

In this study, we tackle musical note estimation for vocal F0 trajectories that tend to have large deviations from original musical scores. The pitches and onset times of musical notes in a musical score can take only discrete values, whereas an F0 trajectory is a continuous signal that can dynamically and smoothly vary over time. F0 trajectories are modulated by vibrato and changes smoothly from one note to another by a portamento. Naive time-frequency quantization of an F0 trajectory therefore outputs a note sequence that often includes statistically-rare chromatic note progressions with unlikely rhythms.

To solve this problem, we propose a statistical method of scale- and rhythm-aware musical note estimation based on integration of a score model describing the process of generating a note sequence and an F0 model describing the process of generating an F0 trajectory from the note sequence (Fig. 1). In the score model, a sequence of musical scales (local keys) is determined by a Markov process and the semitone-level pitch of each note is then determined according to both a scale of the note position and the pitch of a previous note. The onset position of each note on a tatum grid is determined according to that of a previous note to make rhythmic structures. In the F0 model, the time-frequency deviations are added to a step-function-shaped F0 trajectory corresponding to a musical score given by the score model. The integrated model is thus formulated as a hierarchical hidden semi-Markov model (HHSMM). Given a vocal F0 trajectory with a tatum grid, the scales, musical notes, and F0 deviations, which are all latent variables of the proposed model, are jointly estimated by us-



ing a Markov chain Monte Carlo algorithm. A key feature of our method is that musical scales and rhythms work as self-organizing constraints on time-frequency quantization of vocal F0 trajectories.

2. RELATED WORK

In this section, we introduce related work on the analysis of singing voices.

2.1 Vocal F0 Estimation for Music Audio Signals

Estimation of vocal F0 trajectories for music audio signals has actively been studied [1, 3, 7, 9, 14], and the outputs of these methods can be used as inputs of our method. One of the most basic method is subharmonic summation (SHS) [7] that calculates the sum of the harmonic components of each candidate F0. Ikemiya *et al.* [9] improved F0 estimation based on SHS and singing voice separation based on robust principle component analysis (RPCA) [8] by using the mutual dependency of those two tasks. Salamon *et al.* [21] estimated contours of the melody F0 candidates by calculating a salience function and then recursively removed contours which do not form a melody line by using the characteristics of each contour. Durrieu *et al.* [3] extracted a main melody by representing accompaniments with a model inspired by non-negative matrix factorization (NMF) and leading voices with a source-filter model. Mauch *et al.* [14] modified the YIN [1] in a probabilistic way so that the modified system could determine multiple candidate fundamental frequencies and then select one at each frame by using an HMM.

2.2 Musical Note Estimation for Singing Voices

Estimation of musical notes from sung melody has been a hot research topic [6, 11, 13, 15, 17, 18, 20, 23]. A naive method is to take the majority of vocal F0s in each interval of a regular grid [6]. Paiva *et al.* [17] proposed a cascading method based on multipitch detection, multipitch trajectory construction, segmentation of multipitch trajectory, elimination of irrelevant notes, and extraction of notes that form a main melody. Raphael [18] proposed an HMM-based method that estimates pitches, rhythms, and tempos when the number of notes is given. The rhythm and onset deviation models used in [18] are similar to those used in our method. Laaksonen *et al.* [11] divided audio data into segments corresponding to scales and notes by focusing on the boundaries of chords given as input, and independently estimated the notes based on a score function. Ryyänänen *et al.* [20] proposed a method based on a hierarchical HMM in order to capture the different kinds of vocal fluctuations (e.g., vibrato and portamento) within one note. In this model, the transition between pitches is represented in the upper-level HMM and the transition between the vocal fluctuations is represented in the lower-level HMM. Molina *et al.* [15] focused on the hysteresis characteristics of vocal F0s. Nishikimi *et al.* [16] proposed a method based on an HMM that represents the generative process of a vocal F0 trajectory considering the time and frequency deviations. Yang *et al.* [23] proposed a method based on a hierarchical HMM that represents the generative process

of the f_0 - Δf_0 plane. Mauch *et al.* [13] developed a software tool called Tony for extracting pitches. In this tool, a vocal F0 trajectory is estimated by PYIN [14], and musical notes are estimated by a modified version of Ryyänänen's method [20].

3. PROPOSED METHOD

This section explains the proposed method of estimating a sequence of musical notes from a vocal F0 trajectory. The method is based on an HHSMM (Fig. 1) that stochastically generates the F0 trajectory with time-frequency deviations from a sequence of musical notes depending on musical scales. The upper part of the proposed model is an HMM that stochastically generates a sequence of musical notes according to the scales that are assigned to bars. The lower part is an HSMM that represents the musical notes and temporal deviations as latent variables and the frequency deviations as F0 emission probabilities.

3.1 Problem Specification

The problem we tackle is defined as follows:

Input: A vocal F0 trajectory $\mathbf{X} = \{x_t\}_{t=1}^T$ and 16th-note-level tatum $\mathbf{Y} = \{(u_n, v_n)\}_{n=0}^N$,

Output: A sequence of notes $\mathbf{Z} = \{z_j = (p_j, l_j)\}_{j=0}^J$,

where T is the number of frames in a vocal F0 trajectory, x_t is a log frequency at time t , and N is the number of 16th-note-level tatum. $u_n \in \{1, \dots, T+1\}$ is the time of tatum n and the beginning and end of music are represented as $u_0 = 1$ and $u_N = T+1$, respectively. $v_n \in \{0, \dots, 15\}$ is the relative position of tatum n in a bar. J is the number of musical notes estimated by proposed methods, and the j -th note z_j is represented as a pair consisting of an pitch $p_j \in \{1, \dots, K\}$ and a note length $l_j \in \{1, \dots, L\}$ in the unit of tatum, where K is the number of kinds of semitone-level pitches, and p_j indicates any one in $\{\mu_1, \dots, \mu_K\}$, which is a set of log frequencies corresponding to semitone-level pitches. For convenience we introduce the initial note z_0 that does not appear in the actual score.

3.2 Probabilistic Modeling of Musical Scores

This section describes the score model constructed with an HMM that represents rhythms and pitches of musical notes under musical scales.

3.2.1 Modeling Scale Transitions

Scales are represented as $\mathbf{S} = \{s_m\}_{m=0}^M$, where M denotes the number of bars in the musical piece and s_m denotes the scale at the m -th bar. For convenience, we introduce the initial bar s_0 to which the initial note z_0 belongs. Instead of fixing one scale for the whole piece, the scale is allowed to change at bar lines. Each scale s_m takes one of the 24 values of $\{C, C\#, \dots, B\} \times \{\text{major}, \text{minor}\}$. The latent variables \mathbf{S} are described by a Markov chain as

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

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

where $\boldsymbol{\pi} \in \mathbb{R}_{\geq 0}^{24}$ is a set of initial probabilities and $\boldsymbol{\xi}_s \in \mathbb{R}_{\geq 0}^{24}$ is a set of transition probabilities.

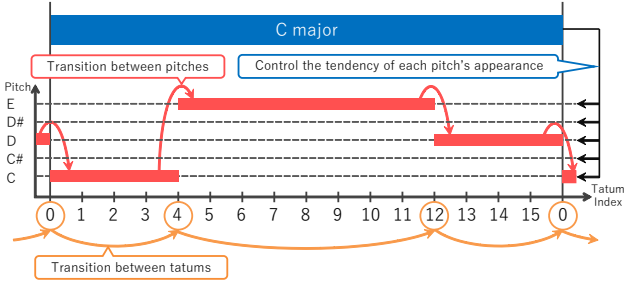


Figure 2: Overview of the score model.

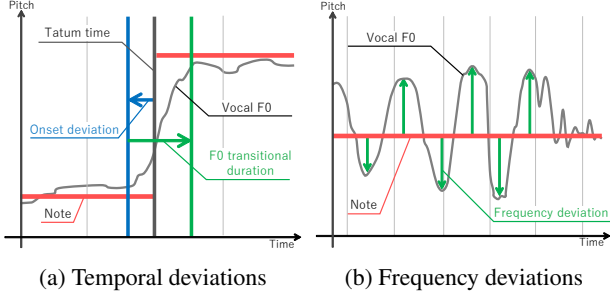


Figure 3: Deviations in a vocal F0 trajectory.

3.2.2 Modeling Pitch Transitions

The sequence of pitches \mathbf{P} is generated by a Markov chain depending on scales \mathbf{S} as follows (Fig. 2):

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

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

where $\phi_s \in \mathbb{R}_{\geq 0}^K$ is a set of initial probabilities, $\psi_{sp} \in \mathbb{R}_{\geq 0}^K$ is a set of transition probabilities, and m is the index of a bar to which the note z_j belongs. Moreover, $\phi_{s_0 p_0}$ and $\psi_{s_m p_{j-1} p_j}$ are defined as

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

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

where $\hat{s} \in \{\text{major}, \text{minor}\}$ is the mode of scale s and $\deg(p; s) \in \{0, \dots, 11\}$ is the degree of pitch p in scale s (defined as the relative pitch class of p from the tonic of scale s). $\hat{\phi}_*$ and $\hat{\psi}_*$ are the initial and transition probabilities of pitch classes, given the scales.

3.2.3 Modeling Onset Transitions

Considering the transition between onset positions of adjacent notes, the model makes \mathbf{Z} have the plausible rhythm. Let $r_{j-1} \in \{v_n\}_{n=1}^N$ be the onset position of the j -th note z_j . The transition probability is given by

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

where the distance between r_{j-1} and r_j indicates the note value l_j of note z_j . We assume that $r_0 = v_0$ and $r_J = v_N$.

3.3 Probabilistic Modeling of F0 Trajectories

The section describes the F0 model based on an HSMM that represents the generative process of a vocal F0 trajectory. In our model, the pitches, onsets, and temporal deviations are represented as latent variables, and the frequency deviations are represented as emission probabilities.

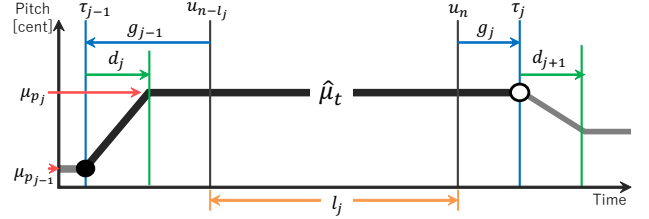


Figure 4: The black bold line represents a sequence of the location parameters of the Cauchy distributions.

3.3.1 Modeling Temporal Deviations

We assume that vocal F0 trajectories include the following two types of temporal deviations (Fig. 3a):

Onset deviation: the gap between the vocal onset time and the note onset time.

F0 transitional duration: the time it takes for singing voices to finish transitioning from one pitch to the next.

The onset deviations $\mathbf{G} = \{g_j\}_{j=0}^J$ accompanying with \mathbf{Z} are represented as discrete latent variables. Each g_j can take an integer value between $-G$ and G . As with the onset position model, g_{j-1} denotes the onset deviation of note z_j . We assume that each g_j is independently generated by

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

where $\boldsymbol{\rho} \in \mathbb{R}_{\geq 0}^{2G+1}$ is a set of onset deviation probabilities. We assume that there are no deviations for the onset of the first note and the offset of the last note, i.e., $g_0 = g_J = 0$.

The F0 transitional durations $\mathbf{D} = \{d_j\}_{j=1}^J$ accompanying with \mathbf{Z} are also represented as discrete latent variables. Each d_j can take a value from 1 to D . The continuous transition of vocal F0s between notes z_{j-1} and z_j is represented by a slanted line spanning d_j frames. We assume that each d_j is independently generated as follows:

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

where $\boldsymbol{\eta} \in \mathbb{R}_{\geq 0}^D$ is a set of duration probabilities.

3.3.2 Modeling Frequency Deviations

The vocal F0 trajectory $\mathbf{X} = \{x_t\}_{t=1}^T$ is generated by imparting probabilistic frequency deviations to the sequence of notes to which probabilistic temporal deviations have already been imparted (Fig. 3b). Assuming that x_t is independently generated at each frame, the emission probability of the j -th note z_j is given by

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

where $x_{\tau' : \tau-1}$ indicates $x_{\tau'}, \dots, x_{\tau-1}$, λ is a scale parameter that represents the scale of the frequency deviations, δ is Kronecker's delta, and $\hat{\mu}_t$ (Fig. 4) is a location parameter given by

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

When the onset of note z_{j+1} is located at the n -th tatum, $\tau_j = u_n + g_j$ and $\tau_{j-1} = u_{n-l_j} + g_{j-1}$.

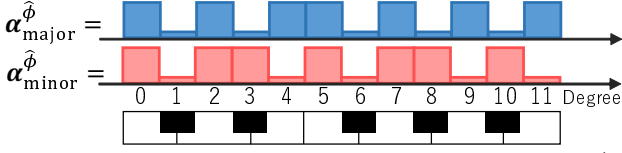


Figure 5: The configuration of the hyperparameter $\mathbf{a}_s^{\hat{\phi}}$.

3.4 Prior Distributions

We put conjugate Dirichlet priors on categorical model parameters π , ξ , $\hat{\phi}$, $\hat{\psi}$, ζ , ρ , and η as follows:

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

where $\mathbf{a}^\pi \in \mathbb{R}_+^{26}$, $\mathbf{a}_s^\xi \in \mathbb{R}_+^{26}$, $\mathbf{a}_s^{\hat{\phi}} \in \mathbb{R}_+^{12}$, $\mathbf{a}_{\text{sdeg}(p;s)}^{\hat{\psi}} \in \mathbb{R}_+^{12}$, $\mathbf{a}_r^\zeta \in \mathbb{R}_+^{16}$, $\mathbf{a}^\rho \in \mathbb{R}_+^{2G+1}$, and $\mathbf{a}^\eta \in \mathbb{R}_+^D$ are hyperparameters. The probability distribution over the 12 pitch classes under a scale is estimated using the priors on the initial and transitional probabilities of those classes. As illustrated in Fig. 5, we set the hyperparameters $\mathbf{a}_s^{\hat{\phi}}$ and $\mathbf{a}_{\text{sdeg}(p;s)}^{\hat{\psi}}$ so that the probability distributions represent the diatonic scales, respectively. Since the Cauchy distribution does not have a conjugate prior, we put a Gamma prior on λ as

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

where a_0^λ and a_1^λ are shape and rate hyperparameters.

3.5 Bayesian Inference

Given an F0 trajectory \mathbf{X} , we aim to calculate the posterior distribution $p(\mathbf{Q}, \mathbf{S}, \Theta | \mathbf{X})$, where $\mathbf{Q} = \{\mathbf{P}, \mathbf{L}, \mathbf{G}, \mathbf{D}\}$ (latent variables) and $\Theta = \{\pi, \xi, \hat{\phi}, \hat{\psi}, \zeta, \rho, \eta\}$ (model parameters). Since this calculation is analytically intractable, we use Markov chain Monte Carlo (MCMC) methods. To get samples of the latent variables \mathbf{S} and \mathbf{Q} , forward filtering-backward sampling algorithms are used. To get samples of Θ except for λ , a set of parameters with conjugate priors, a Gibbs sampling algorithm is used. Since there is no conjugate prior for the parameter λ , we use the Metropolis-Hastings (MH) algorithm. Since \mathbf{S} and \mathbf{Q} share the sequence of notes \mathbf{Z} and are mutually dependent, each variable is updated as follows:

1. Initialize notes \mathbf{Z} with a majority-vote method.
2. Update the sequence of scales \mathbf{S} based on given \mathbf{Z} .
3. Update \mathbf{Q} based on given \mathbf{S} .
4. Update the model parameters Θ .
5. Return to 2.

3.5.1 Inferring Latent Variables \mathbf{S}

Given the sequence of notes \mathbf{Z} , each s_m is sampled in accordance with the probability given by

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

where $s_{m+1:M}$ represents s_{m+1}, \dots, s_M . The calculation of Eq. (14) and sampling of scales \mathbf{S} are performed by the forward filtering-backward sampling method.

In forward filtering, we recursively calculate the probability $\alpha_{s_m}^S$ as follows:

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

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

where j_m is the index of the first note whose onset belongs to the m -th bar. j_m can be calculated from given note values \mathbf{L} .

In backward sampling, Eq. (14) is calculated by using the values calculated in forward filtering, and scales are sampled recursively as follows:

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

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

3.5.2 Inferring Latent Variables \mathbf{Q}

The latent variables \mathbf{Q} can be estimated in a way similar to that in which the latent variables \mathbf{S} are inferred. In forward filtering, we recursively calculate the probability $\alpha_{p_n l_n g_n d_n}^Q$ as follows:

$$\begin{aligned} \alpha_{p_0 l_0 g_0 d_0}^Q &= p(p_0 | \mathbf{S}) = \phi_{y_0 p_0}, \quad (19) \\ \alpha_{p_n l_n g_n d_n}^Q &= p(x_{1:\tau_n-1}, p_n, l_n, g_n, d_n | \mathbf{S}) \\ &= \begin{cases} 0 & (l_n > n) \\ \rho_{g_n} \eta_{d_n} \zeta_{r_0 r_n} \\ \cdot \sum_{p_0} \psi_{s_1 p_0 p_n} e_{p_0 p_n l_n 0 g_n d_n} \alpha_{p_0 l_0 g_0 d_0}^Q & (l_n = n) \\ \sum_{p_{n'}:g_{n'}} \sum_{l_{n'}} \sum_{d_{n'}} \rho_{g_n} \eta_{d_n} \zeta_{r_{n'} r_n} \psi_{s_{m(n')} p_{n'} p_n} \\ \cdot e_{p_{n'} p_n l_n g_{n'} g_n d_n} \alpha_{p_{n'} l_{n'} g_{n'} d_{n'}}^Q & (l_n < n) \end{cases}, \end{aligned} \quad (20)$$

where $\tau_n = u_n + g_n$, $n' = n - l_n$, and $m(n')$ is the index of the bar that the n' -th tatum belongs to. p_n , l_n , g_n , and d_n are the variables of forward messages that correspond to the note whose offset position is at the n -th tatum u_n . Note that these variables are different from j -indexed variables p_j , l_j , g_j , and d_j . Since the onset and offset times of the note $z_n = (p_n, l_n)$ are respectively the $(n-l_n)$ -th tatum and the n -th tatum, the probability $p(l_n)$ which appears in the recursive calculation of Eq. (20) is replaced by $p(r_n | r_{n-l_n})$.

In backward sampling, the posterior distribution of the latent variables is calculated by using the values calculated in forward filtering, and notes and temporal deviations are sampled recursively as follows:

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

3.5.3 Learning Model Parameters Θ

The posterior distributions of the model parameters with the prior distributions are calculated using \mathbf{S} and \mathbf{Q} obtained in the backward sampling steps, and these parameters are sampled according to the posterior distributions as follows:

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

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

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

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

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

where $\mathbf{b}^{\pi} \in \mathbb{R}_{\geq 0}^{26}$ is a unit vector whose s_0 -th element is 1. $\mathbf{b}_s^{\xi} \in \mathbb{R}_{\geq 0}^{26}$ is a vector whose s' -th element indicates the number of transitions between adjacent scales s and s' in the sequence of latent variables \mathbf{Y} . $\mathbf{b}^{\rho} \in \mathbb{R}_{\geq 0}^{2G+1}$ is a vector whose g -th element indicates the number of vocal onset deviations of g in sampled \mathbf{Q} , and $\mathbf{b}^{\eta} \in \mathbb{R}_{\geq 0}^D$ is a vector whose d -th element represents the number of F0 transitional durations of d in sampled \mathbf{Q} . $\mathbf{b}_r^{\zeta} \in \mathbb{R}_{\geq 0}^{16}$ is a vector whose r' -th element represents the number of transitions between adjacent note onset positions r and r' in $\mathbf{R} = \{r_j\}_{j=0}^J$ that can be calculated from the note values \mathbf{L} sampled in backward sampling. Regarding the vector $\mathbf{b}_s^{\hat{\phi}} \in \mathbb{R}_{\geq 0}^{12}$, when the scale of the initial bar and the pitch of the initial note are $s_0 = s$ and $p_0 = p$, the value of the element $b_{s \text{deg}(p;s)}^{\hat{\phi}}$ is 1, and the other elements are 0. Regarding the vector $\mathbf{b}_{s \text{deg}(p;s)}^{\hat{\psi}} \in \mathbb{R}_{\geq 0}^{12}$, the value of $b_{s \text{deg}(p;s) \text{deg}(p';s)}^{\hat{\psi}}$ is increased by one when there is a transition from a pitch p to a pitch p' under a scale s in the sampled latent variables.

To apply the MH sampling to the parameter λ , we define a random-walk proposal distribution as follows:

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

where λ^* is a proposal, λ is the current sample, and γ is a hyperparameter. The proposal λ^* is accepted as the next sample according to the probability given by

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

where $L(\lambda)$ is the complete joint likelihood of λ given by

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

$\{p_j, l_j, g_j, d_j\}_{j=0}^J$ are the values sampled in the backward sampling. The value of λ is updated by λ^* only when the value of $A(\lambda^*, \lambda)$ is larger than a random number sampled from the uniform distribution $\mathcal{U}(0, 1)$.

3.6 Viterbi Decoding

The sequence of latent variables \mathbf{S} and \mathbf{Q} are estimated with the Viterbi algorithm with the model parameters that maximize the joint distribution $p(\mathbf{X}, \mathbf{Q}, \mathbf{S}, \Theta|\Phi)$ in the learning process. As in the inference of latent variables, we initialize \mathbf{Z} by the majority-vote method, \mathbf{S} is estimated based on \mathbf{Z} , and then \mathbf{Q} is estimated depending on the \mathbf{S} estimated in the previous step.

In the Viterbi decoding on scales \mathbf{S} , the value ω_s^S is recursively calculated as follows:

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

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

In the recursive calculation of ω_s^S , the previous state s_{m-1}

that maximizes the value of $\omega_{s_m}^S$ is memorized as $c_{s_m}^S$, and the scales \mathbf{S} are recursively estimated as follows:

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

In the Viterbi decoding on variables \mathbf{Q} , the value ω_{plgd}^Q is recursively calculated as follows:

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

$$\omega_{p_n l_n g_n d_n}^Q = \begin{cases} -\inf & (l_n > n) \\ w^\rho \ln \rho_{g_n} + w^\eta \ln \eta_{d_n} + w^\zeta \ln \zeta_{r_n r_0} \\ \quad + \max_{p_0} \left\{ w^\psi \ln \psi_{s_1 p_0 p_n} \right. \\ \quad \left. + w^e \ln e_{p_0 p_n l_n 0 g_n d_n} + \omega_{p_0 l_0 g_0 d_0}^Q \right\} & (l_n = n), \\ w^\rho \ln \rho_{g_n} + w^\eta \ln \eta_{d_n} + w^\zeta \ln \zeta_{r_n r_{n'}} \\ \quad + \max_{(p_{n'}, l_{n'}, g_{n'}, d_{n'})} \left\{ w^\psi \ln \psi_{s_m(n') p_{n'} p_n} \right. \\ \quad \left. + w^e \ln e_{p_{n'} p_n l_n g_{n'} d_n} + \omega_{p_{n'} l_{n'} g_{n'} d_{n'}}^Q \right\} & (l_n < n) \end{cases} \quad (34)$$

where w^ϕ , w^ψ , w^ρ , w^η , w^ζ , and w^e are the weight parameters that control the balance between probabilities. In the recursive calculation of ω_{plgd}^Q , the previous states $p_{n'}$, $l_{n'}$, $g_{n'}$, and $d_{n'}$ which maximize the value of $\omega_{p_n l_n g_n d_n}^Q$ are memorized as $c_{p_n l_n g_n d_n}^Q$, and the variables \mathbf{Q} are recursively estimated as follows:

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

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

4. EVALUATION

We report comparative experiments conducted to evaluate the performance of the proposed method in musical note estimation from vocal F0 trajectories.

4.1 Experimental Conditions

Among the 100 pieces of popular music in the RWC music database [5], we used 63 pieces that do not include 32nd notes, triplets, harmonizing parts, and overlaps of adjacent notes, which the proposed method cannot deal with. The input F0 trajectories were obtained from the annotation data [4] or automatically estimated by using the state-of-the-art melody extraction method proposed in [9]. The annotation data contain unvoiced regions and the estimation data do not. The tatum times and onset positions were obtained from the annotation data.

The Bayesian inference and Viterbi decoding were independently conducted for each song. The onset transition probabilities were learned in advance from a corpus of rock music [2] without Bayesian learning. The hyperparameters were $\mathbf{a}^\pi = \mathbf{1}$, $\mathbf{a}_s^\xi = \mathbf{1}$, $\mathbf{a}_r^\zeta = \mathbf{1}$, $\mathbf{a}^\rho = \mathbf{a}^\eta = a_0^\lambda = a_1^\lambda = \gamma = 1$, where $\mathbf{1}$ and $\mathbf{1}$ respectively represent the matrix and vector whose elements are all ones. The elements of $\mathbf{a}_s^{\hat{\phi}}$ and $\mathbf{a}_{s \text{deg}(p;s)}^{\hat{\psi}}$ corresponding to musical notes on the scale of \hat{s} were 10 and the others were 1. The weight parameters of the Viterbi algorithm were empirically set as $w^\phi = w^\psi = 29.4$, $w^\rho = 2.4$, $w^\eta = 2.9$, $w^\zeta = 48.5$, and $w^e = 3.8$. To obtain musically-consistent sequences of

Model	Input F0s	Tatum level	Note level
Proposed method	Ground-truth	72.4 ± 1.7	28.1 ± 2.1
	Estimated	68.7 ± 1.3	30.7 ± 1.8
With only rhythms	Ground-truth	71.5 ± 1.6	26.3 ± 2.1
	Estimated	67.7 ± 1.3	29.1 ± 1.8
With only scales	Ground-truth	67.8 ± 1.6	10.6 ± 1.2
	Estimated	65.6 ± 1.2	13.8 ± 1.1
Without scales & rhythms	Ground-truth	67.2 ± 1.5	9.8 ± 1.2
	Estimated	64.6 ± 1.2	12.9 ± 1.1
Majority vote	Ground-truth	54.1 ± 1.5	20.1 ± 1.4
	Estimated	61.0 ± 1.4	22.0 ± 1.5
HMM [16]	Estimated	68.0 ± 1.2	14.8 ± 1.3

Table 1: Average matching rates [%] and their standard errors in tatum and note levels.

musical notes, we put more emphasis on the score model than the F0 model.

For comparison, we tested the majority-vote method as a baseline and the latest conventional method based on a semi-beat-synchronous HMM [16]. Since the conventional method cannot deal with unvoiced regions in a vocal F0 trajectory given as input, we only tested the method for the estimation data. To evaluate the effectiveness of the score model, we tested four versions of the proposed method; a method that does not consider scales (scale transition probabilities) and rhythms (onset transition probabilities), a method considering only scales, a method considering only rhythms, the full method considering both scales and rhythms. To accelerate the inference, the search range of pitches was limited around the pitches estimated by the majority-vote method.

To evaluate the performance of each method, we calculated tatum-level and note-level matching rates by comparing the estimated sequences of musical notes with the ground-truth data. The tatum-level matching rate is the rate of the number of tatum units whose pitches were estimated correctly to the total number of tatum units whose pitches exist in the ground-truth scores. The note-level matching rate is the rate of the number of musical notes whose pitches, onsets, and offsets were estimated correctly to the total number of musical notes in the ground-truth scores. If adjacent notes in the ground-truth scores have the same pitch or are connected by a tie, those notes were regarded as a single note. Since the compared method [16] outputs a pitch in a 16th-note-wise manner, a sequence of the same pitches was regarded as a single note.

4.2 Experimental Results

The experimental results are shown in Table 1¹. The proposed method outperformed the majority-vote method and the conventional method in terms of both measures. Comparing the tatum-level matching rates obtained by the four versions of the proposed method, we confirmed that the score model improved the performance of musical note estimation. The use of the onset transition probabilities

¹ The results of music note estimation by the proposed method are available online: <http://sap.ist.i.kyoto-u.ac.jp/members/nishikimi/demo/ismir2017/>

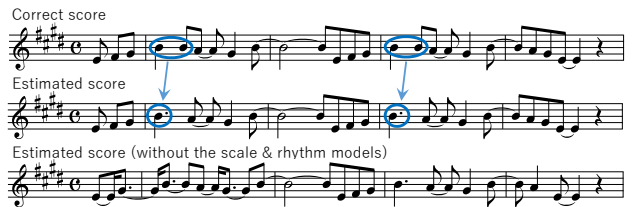


Figure 6: Musical scores estimated from a ground-truth F0 trajectory by the proposed method and its variant without scale and rhythm constraints.

(rhythm constraints) was found to be more effective than that of the scale transition probabilities (scale constraints). Although the tatum-level matching rate obtained by the proposed method (68.7%) was close to that obtained by the conventional method (68.0%), the note-level matching rate obtained by the proposed method (30.7%) was better than that obtained by the conventional method (14.8%). This is a remarkable advantage of the proposed HHSMM that can directly represent both the pitches and durations (onsets and offsets) of musical notes on symbolic musical scores, not on continuous-time piano rolls.

Examples of estimated musical scores are illustrated in Fig. 6. The proposed method yielded the almost accurate musical score except that some notes were merged. To correctly recognize two adjacent notes with the same pitch, it is necessary to refer to original singing voices or music audio signals. The score estimated without considering the score model, on the other hand, included a lot of wrong notes that were inconsistent with music theory. This result also shows the effectiveness of using the score model as musical constraints on musical note estimation.

5. CONCLUSION

This paper presented a statistical method for musical note estimation from a vocal F0 trajectory. Our method is based on an HHSMM that combines a score model (HMM) representing the generative process of a musical score from musical scales with an F0 model (HSMM) representing the generative process of a vocal F0 trajectory with time-frequency deviation from the musical score. We confirmed that the proposed method can yield more musically-consistent sequences of musical notes.

One of the most interesting directions of this research is to use the proposed model as a musically-meaningful prior distribution on a vocal F0 trajectory in vocal F0 estimation for music audio signals. We plan to integrate the proposed “language” model that generates an F0 trajectory from a musical score with an acoustic model that generates a spectrogram from the F0 trajectory in a hierarchical Bayesian manner. This enables us to jointly learn the vocal F0 trajectory and musical score from music audio signals. Joint estimation of beat times and F0s is worth investigating to overcome the problem of estimation-error accumulation in the cascaded estimation approach.

Acknowledgement: This study was partially supported by JSPS KAKENHI Grant Numbers 26700020, 16H01744, and 16J05486 and JST ACCEL No. JPMJAC1602.

6. REFERENCES

- [1] A. de Cheveigné and H. Kawahara. YIN, a fundamental frequency estimator for speech and music. *The Journal of the Acoustical Society of America*, 111(4):1917–1930, 2002.
- [2] T. De Clercq and D. Temperley. A corpus analysis of rock harmony. *Popular Music*, 30(01):47–70, 2011.
- [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*, 18(3):564–575, 2010.
- [4] M. Goto. Aist annotation for the RWC music database. In *The 7th International Conference on Music Information Retrieval (ISMIR 2006)*, pages 359–360, 2006.
- [5] M. Goto, H. Hashiguchi, T. Nishimura, and R. Oka. RWC music database: Popular, classical and jazz music databases. In *The 3rd International Conference on Music Information Retrieval (ISMIR 2002)*, pages 287–288, 2002.
- [6] M. Goto, K. Yoshii, H. Fujihara, M. Mauch, and T. Nakano. Songle: A web service for active music listening improved by user contributions. In *Proc. of the 12th International Society for Music Information Retrieval Conference (ISMIR 2011)*, pages 311–316, 2011.
- [7] Dik J. Hermes. Measurement of pitch by subharmonic summation. *The journal of the acoustical society of America*, 83(1):257–264, 1988.
- [8] P.-S. Huang, S. D. Chen, P. Smaragdis, and M. Hasegawa-Johnson. Singing-voice separation from monaural recordings using robust principal component analysis. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2012)*, pages 57–60, 2012.
- [9] Y. Ikemiya, K. Yoshii, and K. Itoyama. Singing voice analysis and editing based on mutually dependent F0 estimation and source separation. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2015)*, pages 574–578, 2015.
- [10] Y. E. Kim and B. Whitman. Singer identification in popular music recordings using voice coding features. In *3rd International Conference on Music Information Retrieval (ISMIR 2002)*, volume 13, page 17, 2002.
- [11] A. Laaksonen. Automatic melody transcription based on chord transcription. In *Proc. of the 15th International Society for Music Information Retrieval (ISMIR 2014)*, pages 119–124, 2014.
- [12] Y. Li and D. Wang. Separation of singing voice from music accompaniment for monaural recordings. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(4):1475–1487, 2007.
- [13] 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. In *Proc. of the 1st International Conference on Technologies for Music Notation and Representation (TENOR 2015)*, pages 23–30, 2015.
- [14] M. Mauch and S. Dixon. pYIN: A fundamental frequency estimator using probabilistic threshold distributions. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2014)*, pages 659–663, 2014.
- [15] E. Molina, L. J. Tardón, A. M. Barbancho, and I. Barbancho. Siphth: Singing transcription based on hysteresis defined on the pitch-time curve. *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)*, 23(2):252–263, 2015.
- [16] 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. In *Proc. of the 17th International Society for Music Information Retrieval Conference (ISMIR 2016)*, pages 461–467, 2016.
- [17] R. P. Paiva, T. Mendes, and A. Cardoso. On the detection of melody notes in polyphonic audio. In *6th International Conference on Music Information Retrieval (ISMIR 2005)*, pages 175–182, 2005.
- [18] C. Raphael. A graphical model for recognizing sung melodies. In *6th International Conference on Music Information Retrieval (ISMIR 2005)*, pages 658–663, 2005.
- [19] M. Ryyänen, T. Virtanen, J. Paulus, and A. Klauri. Accompaniment separation and karaoke application based on automatic melody transcription. In *2008 IEEE International Conference on Multimedia and Expo*, pages 1417–1420, 2008.
- [20] M. P. Ryyänen and A. P. Klauri. Automatic transcription of melody, bass line, and chords in polyphonic music. *Computer Music Journal*, 32(3):72–86, 2008.
- [21] J. Salamon and E. Gómez. Melody extraction from polyphonic music signals using pitch contour characteristics. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(6):1759–1770, 2012.
- [22] 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*, 14(1):330–341, 2006.
- [23] L. Yang, A. Maezawa, J. B. L. Smith, and E. Chew. Probabilistic transcription of sung melody using a pitch dynamic model. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2017)*, pages 301–305, 2017.