

Generative Modeling of Voice Fundamental Frequency Contours

Hirokazu Kameoka, Kota Yoshizato, Tatsuma Ishihara, Kento Kadowaki, Yasunori Ohishi, and Kunio Kashino

Abstract—This paper introduces a generative model of voice fundamental frequency (F_0) contours that allows us to extract prosodic features from raw speech data. The present F_0 contour model is formulated by translating the Fujisaki model, a well-founded mathematical model representing the control mechanism of vocal fold vibration, into a probabilistic model described as a discrete-time stochastic process. There are two motivations behind this formulation. One is to derive a general parameter estimation framework for the Fujisaki model that allows the introduction of powerful statistical methods. The other is to construct an automatically trainable version of the Fujisaki model that we can incorporate into statistical-model-based text-to-speech synthesizers in such a way that the Fujisaki-model parameters can be learned from a speech corpus in a unified manner. It could also be useful for other speech applications such as emotion recognition, speaker identification, speech conversion and dialogue systems, in which prosodic information plays a significant role. We quantitatively evaluated the performance of the proposed Fujisaki model parameter extractor using real speech data. Experimental results revealed that our method was superior to a state-of-the-art Fujisaki model parameter extractor.

Index Terms—Expectation-maximization algorithm, Fujisaki model, prosody, voice fundamental frequency contour.

I. INTRODUCTION

PROSODY assists the listener to interpret an utterance by grouping words into larger information units and drawing attention to specific words. It also plays an important role in conveying various types of non-linguistic information such as the identity, intention, attitude and mood of the speaker. Since the voice fundamental frequency (F_0) contour is an important

acoustic correlate of many prosodic constructs, modeling and analyzing F_0 contours is potentially useful for many speech applications such as speech synthesis, speaker identification, speech conversion and dialogue systems, in which prosodic information plays a significant role. It is also important to note that F_0 contours indicate intonation in pitch accent languages.

An F_0 contour consists of long-term pitch variations over the duration of prosodic units and short-term pitch variations in accented syllables. The former usually contribute to phrasing while the latter contribute to accentuation during an utterance. These two types of pitch variations can be interpreted as the manifestations of two independent movements by the thyroid cartilage. The Fujisaki model [5], [6] is a well-founded mathematical model that describes an F_0 contour as the sum of these two contributions. This model approximates actual F_0 contours of speech fairly well when the model parameters are appropriately chosen, and its validity has been demonstrated for many typologically diverse languages [6]–[13]. Since prosodic features in speech are predominantly characterized by the levels and timings of the phrase and accent components, one important challenge is to solve the inverse problem of estimating the Fujisaki-model parameters automatically from a raw F_0 contour.

However, this problem has proved difficult to solve. Several techniques have already been developed [6], [14]–[18], but so far with limited success due to the difficulty of finding optimal parameters under the constraints imposed in the Fujisaki model. While the Fujisaki model describes a deterministic process for generating voice F_0 contours, this paper proposes formulating a stochastic counterpart of the Fujisaki model. As will be shown in the subsequent sections, this makes it possible to use powerful statistical inference techniques for estimating the underlying parameters of the Fujisaki model. Another important motivation for this formulation is to construct an automatically trainable version of the Fujisaki model that we can smoothly incorporate into text-to-speech synthesis systems or speech conversion systems so as to guarantee the naturalness of computer-generated speech.

The rest of this paper is organized as follows. Section II briefly reviews the original Fujisaki model. Section III describes a discrete-time version of the Fujisaki model, on which basis Section IV formulates a generative model of voice F_0 contours. Section V presents two iterative algorithms, which locally maximize the posterior density of the Fujisaki model parameters given an observed F_0 contour. Section VI presents experimental evaluations of the present method in terms of its ability as a Fujisaki model parameter extractor.

Manuscript received October 22, 2013; revised January 09, 2015; accepted March 13, 2015. Date of publication April 01, 2015; date of current version April 15, 2015. This work was supported by the JSPS KAKENHI under Grants 26730100 and 26280060. Earlier versions of this work were presented at SAPA 2010 [1], Speech Prosody 2012 [2], Interspeech 2012 [3], and Interspeech 2013 [4] as workshop/conference papers. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Nobutaka Ono.

H. Kameoka is with the Graduate School of Information Science and Technology, The University of Tokyo, 113-8656, Japan, and also with NTT Communication Science Laboratories, NTT Corporation, 243-0198, Japan (e-mail: kameoka@hil.t.u-tokyo.ac.jp; kameoka.hirokazu@lab.ntt.co.jp).

K. Yoshizato, T. Ishihara, and K. Kadowaki are with the Graduate School of Information Science and Technology, The University of Tokyo, 113-8656, Japan (e-mail: yoshizato@hil.t.u-tokyo.ac.jp; ishihara@hil.t.u-tokyo.ac.jp; kadowaki@hil.t.u-tokyo.ac.jp).

Y. Ohishi and K. Kashino are with the NTT Communication Science Laboratories, NTT Corporation, 243-0198, Japan (e-mail: ohishi.yasunori@lab.ntt.co.jp; kashino.kunio@lab.ntt.co.jp).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TASLP.2015.2418576

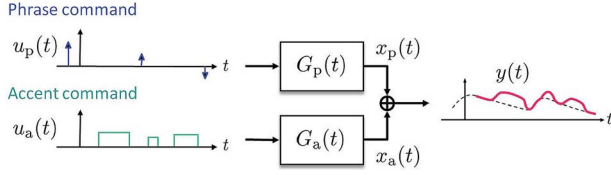


Fig. 1. Original Fujisaki model [5].

II. ORIGINAL FUJISAKI MODEL

The Fujisaki model [5], shown in Fig. 1, assumes that an F_0 contour on a logarithmic scale, $y(t)$, where t is time, is the superposition of three components: a phrase component $x_p(t)$, an accent component $x_a(t)$, and a base value μ_b :

$$y(t) = x_p(t) + x_a(t) + \mu_b. \quad (1)$$

The phrase and accent components are considered to be associated with mutually independent types of movement of the thyroid cartilage with different degrees of freedom and muscular reaction times. The phrase component $x_p(t)$ consists of the large-scale pitch variations over the duration of the prosodic units, and the accent component $x_a(t)$ consists of the smaller-scale pitch variations in accented syllables. These two components are modeled as the outputs of second-order critically damped filters; one being excited with a command function $u_p(t)$ consisting of Dirac deltas (phrase commands), and the other with $u_a(t)$ consisting of rectangular pulses (accent commands):

$$x_p(t) = G_p(t) * u_p(t), \quad (2)$$

$$x_a(t) = G_a(t) * u_a(t), \quad (3)$$

where $*$ denotes convolution over time and

$$G_p(t) = \begin{cases} \alpha^2 t e^{-\alpha t} & (t \geq 0) \\ 0 & (t < 0) \end{cases}, \quad (4)$$

$$G_a(t) = \begin{cases} \beta^2 t e^{-\beta t} & (t \geq 0) \\ 0 & (t < 0) \end{cases}, \quad (5)$$

μ_b is a constant value related to the lower bound of the speaker's F_0 , below which no regular vocal fold vibration can be maintained. α and β are natural angular frequencies of the two second-order systems, which are known to be almost constant within an utterance as well as across utterances for a particular speaker. It has been shown that $\alpha = 3$ rad/s and $\beta = 20$ rad/s can be used as default values [6], [15].

III. DISCRETIZED FUJISAKI MODEL

In this section, we apply a backward difference s -to- z transform to the phrase and accent control mechanisms described as continuous-time linear systems in order to obtain a discrete-time version of the Fujisaki model. The reason for the discretization will be made apparent later. The transfer function (the Laplace transform of the impulse response) of the phrase control mechanism is given in the s -domain as

$$\mathcal{G}_p(s) = \mathcal{L}[G_p(t)] = \frac{\alpha^2}{(s + \alpha)^2}. \quad (6)$$

The backward difference transform approximates the time differential operator s by the backward difference operator in the z -domain such that

$$s \simeq \frac{1 - z^{-1}}{t_0}, \quad (7)$$

where t_0 is the sampling period of the discrete-time representation. By undertaking this transform, the transfer function of the inverse system $\mathcal{G}_p^{-1}(s)$ can be written in the z -domain as

$$\mathcal{G}_p^{-1}(z) = f_{p2}z^{-2} + f_{p1}z^{-1} + f_{p0}, \quad (8)$$

where

$$f_{p2} = (\psi - 1)^2, \quad (9)$$

$$f_{p1} = -2\psi(\psi - 1), \quad (10)$$

$$f_{p0} = \psi^2, \quad (11)$$

$$\psi = 1 + 1/(\alpha t_0). \quad (12)$$

Let us use $u_p[k]$ and $x_p[k]$ to denote the discrete-time version of the phrase command function and phrase component, respectively, where k is the discrete-time index. $x_p[k]$ can thus be regarded as the output of a constrained all-pole system whose characteristics are governed by a single parameter ψ (or α)

$$u_p[k] = f_{p0}x_p[k] + f_{p1}x_p[k-1] + f_{p2}x_p[k-2]. \quad (13)$$

In the same way, the relationship between the accent command function $u_a[k]$ and the accent component $x_a[k]$ is described as

$$u_a[k] = f_{a0}x_a[k] + f_{a1}x_a[k-1] + f_{a2}x_a[k-2], \quad (14)$$

where

$$f_{a2} = (\varphi - 1)^2, \quad (15)$$

$$f_{a1} = -2\varphi(\varphi - 1), \quad (16)$$

$$f_{a0} = \varphi^2, \quad (17)$$

$$\varphi = 1 + 1/(\beta t_0). \quad (18)$$

Altogether, the discrete-time version of the Fujisaki model can be expressed as the superposition of the three components: $x_p[k] + x_a[k] + \mu_b$.

IV. GENERATIVE MODEL OF VOICE F_0 CONTOURS

Here, we model the probabilistic generative process of a speech F_0 contour based on the discrete-time version of the Fujisaki model.

A. Modeling Phrase and Accent Command Pair

We first describe the process for generating the phrase and accent command functions, $u_p[k]$ and $u_a[k]$. In the original Fujisaki model, they must satisfy the following requirements:

- 1) Phrase commands are a set of impulses and accent commands are a set of step-wise functions.
- 2) A phrase command occurs at the start of an utterance or after the offset of an accent command in the preceding phrase, and is followed by the onset of the next accent command. This means that a phrase command will not occur while an accent command is being activated.
- 3) The onset of an accent command is followed by its offset. This means that neighboring accent commands will not overlap each other.

According to assumption 2, $u_p[k]$ and $u_a[k]$ are reciprocally constrained and so they should not simply be modeled separately. One challenge as regards the estimation of the Fujisaki model parameters has been how to deal with the optimization problem under these constraints. As a convenient way of incorporating these requirements into the command functions,

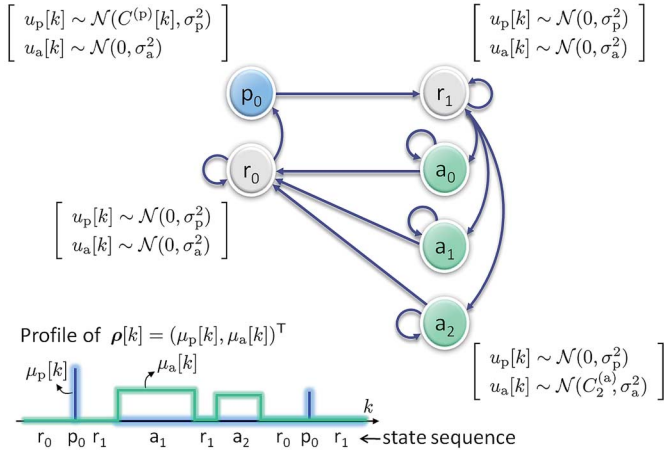


Fig. 2. Command function modeling with HMM.

we propose modeling the $u_p[k]$ and $u_a[k]$ pair using a hidden Markov model (HMM).

We denote the $u_p[k]$ and $u_a[k]$ pair by $\mathbf{o}[k]$ and assume that it is normally distributed:

$$\mathbf{o}[k] \sim \mathcal{N}(\mathbf{o}[k]; \boldsymbol{\rho}[k], \boldsymbol{\Upsilon}), \quad (19)$$

where

$$\mathbf{o}[k] = \begin{bmatrix} u_p[k] \\ u_a[k] \end{bmatrix}, \quad \boldsymbol{\rho}[k] = \begin{bmatrix} \mu_p[k] \\ \mu_a[k] \end{bmatrix}, \quad \boldsymbol{\Upsilon} = \begin{bmatrix} \sigma_p^2 & 0 \\ 0 & \sigma_a^2 \end{bmatrix}.$$

Eq. (19) can be viewed as an HMM in which the output distribution of each state is a Gaussian distribution and the mean vector $\boldsymbol{\rho}[k]$ evolves over time as a result of the state transition. The mean vector $\boldsymbol{\rho}[k]$ consists of the means of the phrase and accent command functions, $\mu_p[k]$ and $\mu_a[k]$. The use of an HMM allows us to incorporate assumptions 1–3 into $\mu_p[k]$ and $\mu_a[k]$ by constraining the path of the state transitions as illustrated in Fig. 2.

The present HMM consists of $N + 3$ distinct states, r_0 , p_0 , r_1 and a_0, \dots, a_{N-1} . N is the number of possible values that the magnitude of each accent command can take. It can thus be understood as the resolution of magnitude “quantization”: the larger it becomes, the more finely the model is able to express the accent command function. In state r_0 , $\mu_p[k]$ and $\mu_a[k]$ are both restricted to zero. In state p_0 , $\mu_p[k]$ can take a non-zero value, $C^{(p)}[k]$, whereas $\mu_a[k]$ is still restricted to zero. In state p_0 , no self-transitions are allowed. In state r_1 , $\mu_p[k]$ and $\mu_a[k]$ become zero again. This path constraint restricts $\mu_p[k]$ to consisting of isolated deltas. State r_1 leads to states a_0, \dots, a_{N-1} , in each of which $\mu_a[k]$ can take a different non-zero value $C_n^{(a)}$, whereas $\mu_p[k]$ is forced to be zero. Direct state transitions from state a_n to state a'_n ($n \neq n'$) without passing through state r_1 are not allowed. This constraint restricts $\mu_a[k]$ to consisting of rectangular pulses. It should also be noted that this HMM ensures that no more than one command will be active at each point in time. The use of the HMM described above for modeling the command functions has been our primary reason for translating the Fujisaki model into its discrete-time counterpart.

The state segments correspond to the timings and durations of phrase and accent commands. If the statistical distributions of the state durations can be trained a priori, they can be useful

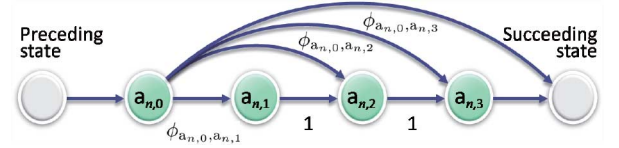


Fig. 3. A duration-explicit representation of the hidden states. Splitting state a_n into substates $a_{n,0}$, $a_{n,1}$, $a_{n,2}$, and $a_{n,3}$ allows us to parametrize the duration of each hidden state. For example, $\phi_{a_{n,0},a_{n,1}}$ corresponds to the probability of staying at state a_n with 4 consecutive times.

for estimating the timings of phrase and accent commands. While an ordinary HMM assumes the state durations to be geometrically distributed, it would be more convenient if we were allowed to assume arbitrary distributions. To allow arbitrary distributions, we propose splitting each state into a certain number of substates such that they all have exactly the same emission densities. Fig. 3 shows an example of the splitting of state a_n . The number of substates is set at a sufficiently large value and the transition probability from substate $a_{n,m}$ to substate $a_{n,m+1}$ is set at 1 for $m \neq 0$. This state splitting allows us to assume arbitrary distributions over the durations for which the process stays in state a_n through the settings of the transition probability. The transition probability from substate $a_{n,0}$ to substate $a_{n,m}$ ($m \geq 1$) corresponds to the probability of the present HMM generating a rectangular pulse that has a particular duration. In the same way, we split states r_0 and r_1 to parameterize the probability of the spacing between phrase and accent commands. Note that this is equivalent to the explicit-duration HMM proposed by Ferguson [19]. Alternatively, the use of a hidden semi-Markov model [20], [21] would also be appropriate for the same purpose. Henceforth, we use the notation $r_0 = \{r_{0,0}, r_{0,1}, \dots\}$, $r_1 = \{r_{1,0}, r_{1,1}, \dots\}$, and $a_n = \{a_{n,0}, a_{n,1}, \dots\}$. Let $\phi_{i',i}$ be the logarithm of the transition probability from state i' and i . To sum up, the present HMM is defined as follows:

Output sequence: $\{\mathbf{o}[k]\}_{k=0}^{K-1}$

Set of states: $\mathcal{S} = \{r_0, p_0, r_1, a_0, \dots, a_{N-1}\}$

State sequence: $\{s_k\}_{k=0}^{K-1}$

Output distribution: $P(\mathbf{o}[k] | s_k = i) = \mathcal{N}(\mathbf{c}_i[k], \boldsymbol{\Upsilon})$

$$\mathbf{c}_i[k] = \begin{cases} (0, 0)^\top & (i \in r_0) \\ (C^{(p)}[k], 0)^\top & (i = p_1) \\ (0, 0)^\top & (i \in r_1) \\ (0, C_n^{(a)})^\top & (i \in a_n) \end{cases}$$

$$\boldsymbol{\Upsilon} = \begin{bmatrix} \sigma_p^2 & 0 \\ 0 & \sigma_a^2 \end{bmatrix}$$

Transition probability: $\phi_{i',i} = \log P(s_k = i | s_{k-1} = i')$

The free parameters to be estimated in our command function model consist of the state sequence, $\{s_k\}_{k=0}^{K-1}$, and the mean and variance of the output distribution of each state, $\{C^{(p)}[k]\}_{k=0}^{K-1}$, $\{C_n^{(a)}\}_{n=0}^{N-1}$, $\{\sigma_p^2, \sigma_a^2\}$. Hereafter, we use \mathbf{s} to denote $\{s_k\}_{k=0}^{K-1}$ and θ to denote the rest of the parameters:

$$\mathbf{s} := \{s_k\}_{k=1}^K,$$

$$\theta := \{\{C^{(p)}[k]\}_{k=0}^{K-1}, \{C_n^{(a)}\}_{n=0}^{N-1}, \sigma_p^2, \sigma_a^2\}.$$

The generating process for the phrase and accent components is summarized as follows: First, the state sequence $\{s_k\}_{k=0}^{K-1}$ is generated according to a Markov chain. Given the state sequence $\{s_k\}_{k=0}^{K-1}$, the mean sequences $\{\mu_p[k]\}_{k=0}^{K-1}$ and $\{\mu_a[k]\}_{k=0}^{K-1}$ are determined by

$$\begin{bmatrix} \mu_p[k] \\ \mu_a[k] \end{bmatrix} = \boldsymbol{\rho}[k] = \mathbf{c}_{s_k}[k]. \quad (20)$$

Given $\boldsymbol{\rho}[k]$ and Υ the present HMM generates the $u_p[k]$ and $u_a[k]$ pair according to Eq. (19). From Eq. (13) and Eq. (14), $u_p[k]$ and $u_a[k]$ are then fed through different all-pole systems to generate the phrase and accent components, $x_p[k]$ and $x_a[k]$.

B. Uncertainty of F_0 Observations

In the following, we assume that F_0 contours and voiced/unvoiced (V/UV) segments are obtained in advance by using a pitch extractor and a V/UV detector. For real speech, F_0 values should not always be considered reliable. For example, F_0 estimates obtained with a pitch extractor in unvoiced regions would be totally unreliable. When performing parameter inference, we should trust only reliable observations and neglect unreliable ones. To incorporate the degree of uncertainty of F_0 observations, we consider modeling an observed F_0 contour $y[k]$ as a superposition of the ‘‘ideal’’ F_0 contour, i.e., $x_p[k] + x_a[k] + \mu_b$, and a normally distributed noise component

$$x_n[k] \sim \mathcal{N}(0, \sigma_n^2[k]), \quad (21)$$

where $\sigma_n^2[k]$ represents the degree of uncertainty of the F_0 observation at time k , which is assumed to be given. For example, one simple way would be to set $\sigma_n^2[k]$ at a small value near 0 for voiced regions and a sufficiently large value for unvoiced regions. By denoting

$$x_b[k] = \mu_b + x_n[k], \quad (22)$$

the entire F_0 contour is given by

$$y[k] = x_p[k] + x_a[k] + x_b[k]. \quad (23)$$

C. Likelihood Function and Prior Probabilities

In this subsection, we derive the probability density function of an observed F_0 contour, $y[0], \dots, y[K-1]$, based on the probabilistic modeling of the command functions and the reliability modeling presented in the previous subsections. From Eq. (19),

$$u_p[k]|\theta, s_k \sim \mathcal{N}(\mu_p[k], \sigma_p^2), \quad (24)$$

$$u_a[k]|\theta, s_k \sim \mathcal{N}(\mu_a[k], \sigma_a^2). \quad (25)$$

If we let \mathbf{u}_p and \mathbf{u}_a be

$$\mathbf{u}_p = (u_p[0], \dots, u_p[K-1])^\top, \quad (26)$$

$$\mathbf{u}_a = (u_a[0], \dots, u_a[K-1])^\top, \quad (27)$$

we can write Eqs. (24) and (25) in vector notation:

$$\mathbf{u}_p|\theta, \mathbf{s} \sim \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p), \quad (28)$$

$$\mathbf{u}_a|\theta, \mathbf{s} \sim \mathcal{N}(\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a), \quad (29)$$

where

$$\boldsymbol{\mu}_p = (\mu_p[0], \dots, \mu_p[K-1])^\top, \quad (30)$$

$$\boldsymbol{\mu}_a = (\mu_a[0], \dots, \mu_a[K-1])^\top, \quad (31)$$

$$\boldsymbol{\Sigma}_p = \sigma_p^2 \mathbf{I}, \quad (32)$$

$$\boldsymbol{\Sigma}_a = \sigma_a^2 \mathbf{I}. \quad (33)$$

By using the linear equation given by Eqs. (13) and (14), the vectors consisting of the phrase and accent components

$$\mathbf{x}_p = (x_p[0], \dots, x_p[K-1])^\top, \quad (34)$$

$$\mathbf{x}_a = (x_a[0], \dots, x_a[K-1])^\top, \quad (35)$$

can be written in terms of \mathbf{u}_p and \mathbf{u}_a ,

$$\mathbf{u}_p = \mathbf{F}_p \mathbf{x}_p, \quad (36)$$

$$\mathbf{u}_a = \mathbf{F}_a \mathbf{x}_a, \quad (37)$$

where

$$\mathbf{F}_p := \begin{bmatrix} f_{p0} & & & & O \\ f_{p1} & f_{p0} & & & \\ f_{p2} & f_{p1} & f_{p0} & & \\ & \ddots & \ddots & \ddots & \\ O & & f_{p2} & f_{p1} & f_{p0} \end{bmatrix}, \quad (38)$$

$$\mathbf{F}_a := \begin{bmatrix} f_{a0} & & & & O \\ f_{a1} & f_{a0} & & & \\ f_{a2} & f_{a1} & f_{a0} & & \\ & \ddots & \ddots & \ddots & \\ O & & f_{a2} & f_{a1} & f_{a0} \end{bmatrix}. \quad (39)$$

Hence, it follows from Eqs. (28), (29), (36) and (37) that

$$\mathbf{x}_p|\theta, \mathbf{s}, \psi \sim \mathcal{N}(\mathbf{F}_p^{-1} \boldsymbol{\mu}_p, \mathbf{F}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{F}_p^{-1})^\top), \quad (40)$$

$$\mathbf{x}_a|\theta, \mathbf{s}, \varphi \sim \mathcal{N}(\mathbf{F}_a^{-1} \boldsymbol{\mu}_a, \mathbf{F}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{F}_a^{-1})^\top). \quad (41)$$

We refer to $x_b[k]$ as the base component and let \mathbf{x}_b be

$$\mathbf{x}_b = (x_b[0], \dots, x_b[K-1])^\top. \quad (42)$$

It follows from Eqs. (21) and (22) that \mathbf{x}_b is normally distributed

$$\mathbf{x}_b|\mu_b \sim \mathcal{N}(\mu_b \mathbf{1}, \boldsymbol{\Sigma}_b), \quad (43)$$

where

$$\mathbf{1} = (1, \dots, 1)^\top, \quad (44)$$

$$\boldsymbol{\Sigma}_b = \text{diag}(\sigma_n^2[0], \dots, \sigma_n^2[K-1]). \quad (45)$$

Let us define a vector consisting of observed F_0 s as

$$\mathbf{y} = (y[0], \dots, y[K-1])^\top. \quad (46)$$

Hence,

$$\mathbf{y} = \mathbf{x}_p + \mathbf{x}_a + \mathbf{x}_b. \quad (47)$$

We can conclude from Eqs. (47), (40), (41) and (43) that \mathbf{y} is normally distributed such that

$$\mathbf{y}|\Theta \sim \mathcal{N}(\mathbf{F}_p^{-1} \boldsymbol{\mu}_p + \mathbf{F}_a^{-1} \boldsymbol{\mu}_a + \mu_b \mathbf{1}, \mathbf{F}_p^{-1} \boldsymbol{\Sigma}_p (\mathbf{F}_p^{-1})^\top + \mathbf{F}_a^{-1} \boldsymbol{\Sigma}_a (\mathbf{F}_a^{-1})^\top + \boldsymbol{\Sigma}_b), \quad (48)$$

where $\Theta := \{\theta, \mathbf{s}, \psi, \varphi, \mu_b\}$. Overall, the likelihood function of the Fujisaki model parameters Θ given \mathbf{y} can be written as

$$P(\mathbf{y}|\Theta) = \frac{1}{\sqrt{(2\pi)^K |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\},$$

$$\boldsymbol{\mu} = \mathbf{F}_p^{-1} \boldsymbol{\mu}_p + \mathbf{F}_a^{-1} \boldsymbol{\mu}_a + \mu_b \mathbf{1},$$

$$\Sigma = \mathbf{F}_p^{-1} \Sigma_p (\mathbf{F}_p^\top)^{-1} + \mathbf{F}_a^{-1} \Sigma_a (\mathbf{F}_a^\top)^{-1} + \Sigma_b. \quad (49)$$

As for the prior probability of Θ , we assume that the phrase control parameter ψ , accent control parameter φ and state sequence $\{s[k]\}_{k=1}^K$ are independent of each other. Recall that we assumed in IV-A that $\{s[k]\}_{k=1}^K$ is a first-order Markov chain. We further assume that all other parameters are uniformly distributed. Thus,

$$P(\Theta) \propto P(\psi)P(\varphi)P(\mathbf{s}), \quad (50)$$

$$P(\mathbf{s}) = P(s_0) \prod_{k=1}^{K-1} P(s_k | s_{k-1}). \quad (51)$$

V. PARAMETER OPTIMIZATION ALGORITHM

A. Expectation-Maximization (EM) approach

Here we describe an iterative algorithm, which locally maximizes the posterior density of Θ given \mathbf{y} , $P(\Theta|\mathbf{y}) \propto P(\mathbf{y}|\Theta)P(\Theta)$. By regarding the set consisting of the phrase, accent and base components, $\mathbf{x} := (\mathbf{x}_p^\top, \mathbf{x}_a^\top, \mathbf{x}_b^\top)^\top$, as the complete data, this problem can be viewed as an incomplete data problem, which can be dealt with using the Expectation-Maximization (EM) algorithm [22], [23].

The log-likelihood of Θ given the complete data is given as

$$\log P(\mathbf{x}|\Theta) \stackrel{\ominus}{=} \frac{1}{2} \log |\Lambda^{-1}| - \frac{1}{2} (\mathbf{x} - \mathbf{m})^\top \Lambda^{-1} (\mathbf{x} - \mathbf{m}),$$

$$\mathbf{x} := \begin{bmatrix} \mathbf{x}_p \\ \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \mathbf{m} := \begin{bmatrix} \mathbf{F}_p^{-1} \boldsymbol{\mu}_p \\ \mathbf{F}_a^{-1} \boldsymbol{\mu}_a \\ \mu_b \mathbf{1} \end{bmatrix},$$

$$\Lambda^{-1} := \begin{bmatrix} \mathbf{F}_p^\top \Sigma_p^{-1} \mathbf{F}_p & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{F}_a^\top \Sigma_a^{-1} \mathbf{F}_a & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \Sigma_b^{-1} \end{bmatrix}, \quad (52)$$

where $\stackrel{\xi}{=}$ denotes equality up to a term independent of ξ . Taking the conditional expectation of Eq. (52) with respect to \mathbf{x} given \mathbf{y} and $\Theta = \Theta'$, and then adding $\log P(\Theta)$ to both sides, we obtain an auxiliary function

$$Q(\Theta, \Theta') \stackrel{\ominus}{=} \frac{1}{2} [\log |\Lambda^{-1}| - \text{tr}(\Lambda^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{y}; \Theta'])]$$

$$+ 2\mathbf{m}^\top \Lambda^{-1} \mathbb{E}[\mathbf{x}\mathbf{y}; \Theta'] - \mathbf{m}^\top \Lambda^{-1} \mathbf{m} + \log P(\Theta). \quad (53)$$

Because the relationship between the incomplete data and the complete data can be written as $\mathbf{y} = \mathbf{H}\mathbf{x}$ where $\mathbf{H} := [\mathbf{I}, \mathbf{I}, \mathbf{I}]$, $\mathbb{E}[\mathbf{x}\mathbf{y}; \Theta]$ and $\mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{y}; \Theta]$ are given explicitly as

$$\mathbb{E}[\mathbf{x}\mathbf{y}; \Theta] = \mathbf{m} + \Lambda \mathbf{H}^\top (\mathbf{H} \Lambda \mathbf{H}^\top)^{-1} (\mathbf{y} - \mathbf{H}\mathbf{m}), \quad (54)$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{y}; \Theta] = \Lambda - \Lambda \mathbf{H}^\top (\mathbf{H} \Lambda \mathbf{H}^\top)^{-1} \mathbf{H} \Lambda$$

$$+ \mathbb{E}[\mathbf{x}\mathbf{y}; \Theta] \mathbb{E}[\mathbf{x}\mathbf{y}; \Theta]^\top. \quad (55)$$

It can be shown that an iterative procedure consisting of maximizing $Q(\Theta, \Theta')$ with respect to Θ (the maximization step) and substituting Θ into Θ' (the expectation step) locally maximizes the posterior density $P(\Theta|\mathbf{y})$. The expectation step computes $\mathbb{E}[\mathbf{x}\mathbf{y}; \Theta']$ and $\mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{y}; \Theta']$ according to Eqs. (54) and (55) by substituting the current parameter estimate into Θ' .

Now, if we partition $\mathbb{E}[\mathbf{x}\mathbf{y}; \Theta']$ into three $K \times 1$ blocks and $\mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{y}; \Theta']$ into nine $K \times K$ blocks such that

$$\mathbb{E}[\mathbf{x}\mathbf{y}; \Theta'] = \begin{bmatrix} \bar{\mathbf{x}}_p \\ \bar{\mathbf{x}}_a \\ \bar{\mathbf{x}}_b \end{bmatrix}, \mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{y}; \Theta'] = \begin{bmatrix} \mathbf{R}_p & * & * \\ * & \mathbf{R}_a & * \\ * & * & \mathbf{R}_b \end{bmatrix}, \quad (56)$$

where $*$ stands for blocks that we can ignore hereafter, the auxiliary function can be rewritten in a more convenient form:

$$Q(\Theta, \Theta') \stackrel{\ominus}{=} \frac{1}{2} [\log |\mathbf{F}_p^\top \Sigma_p^{-1} \mathbf{F}_p| + \log |\mathbf{F}_a^\top \Sigma_a^{-1} \mathbf{F}_a| + \log |\Sigma_b^{-1}|]$$

$$- \text{tr}(\mathbf{F}_p^\top \Sigma_p^{-1} \mathbf{F}_p \mathbf{R}_p) + 2\boldsymbol{\mu}_p^\top \Sigma_p^{-1} \mathbf{F}_p \bar{\mathbf{x}}_p$$

$$- \text{tr}(\mathbf{F}_a^\top \Sigma_a^{-1} \mathbf{F}_a \mathbf{R}_a) + 2\boldsymbol{\mu}_a^\top \Sigma_a^{-1} \mathbf{F}_a \bar{\mathbf{x}}_a$$

$$- \text{tr}(\Sigma_b^{-1} \mathbf{R}_b) + 2\mu_b \mathbf{1}^\top \Sigma_b^{-1} \bar{\mathbf{x}}_b$$

$$- \boldsymbol{\mu}_p^\top \Sigma_p^{-1} \boldsymbol{\mu}_p - \boldsymbol{\mu}_a^\top \Sigma_a^{-1} \boldsymbol{\mu}_a - \mu_b^2 \mathbf{1}^\top \Sigma_b^{-1} \mathbf{1}]$$

$$+ \log P(\Theta). \quad (57)$$

The update formula for each parameter in the maximization step can be derived using Eq. (57).

- 1) **State sequence** s_0, \dots, s_{K-1} : Leaving only the terms in $Q(\Theta, \Theta')$ that depend on $s := \{s_k\}_{k=0}^{K-1}$, we have

$$Q(\Theta, \Theta') \stackrel{s}{=} -\frac{1}{2} \sum_{k=0}^{K-1} (\mathbf{o}[k] - \mathbf{c}_{s_k}[k])^\top \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[k] - \mathbf{c}_{s_k}[k])$$

$$+ \phi_{s_0} + \sum_{k=1}^{K-1} \phi_{s_{k-1}, s_k}, \quad (58)$$

where $\mathbf{o}[k] := ([F_p \bar{\mathbf{x}}_p]_k, [F_a \bar{\mathbf{x}}_a]_k)^\top$. Here the notation $[\cdot]_k$ is used to denote the k -th element of a vector. The state sequence $\{s_k\}_{k=0}^{K-1}$ maximizing $Q(\Theta, \Theta')$ can be solved efficiently using the Viterbi algorithm as follows. We first set $\delta_0(i)$ at

$$\delta_0(i) = -\frac{1}{2} (\mathbf{o}[0] - \mathbf{c}_i[0])^\top \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[0] - \mathbf{c}_i[0]) + \phi_i. \quad (59)$$

for all the hidden states i . If we consider state $r_{0,0}$ to be the starting state, we shall set ϕ_i at

$$\phi_i = \begin{cases} 0 & (i = r_{0,0}) \\ -\infty & (i \neq r_{0,0}) \end{cases}. \quad (60)$$

We can compute $\delta_k(i)$ for $k = 1, \dots, K-1$ recursively via

$$\delta_k(i) = \max_i [\delta_{k-1}(i') + \phi_{i', i}]$$

$$- \frac{1}{2} (\mathbf{o}[k] - \mathbf{c}_i[k])^\top \boldsymbol{\Upsilon}^{-1} (\mathbf{o}[k] - \mathbf{c}_i[k]). \quad (61)$$

The most likely transition for each state should be registered at each recursion $\Psi_k(i) = \arg \max_{i'} [\delta_{k-1}(i') + \phi_{i', i}]$, so that the most likely state sequence can be traced at the end of the recursion with $s_{k-1} = \Psi_k(s_k)$ ($k = K-1, \dots, 1$), where $s_K = \arg \max_i \delta_K(i)$. Substituting

the updated state sequence $\{s_k\}$ into Eq. (20), we finally obtain the updated μ_p and μ_a .

- 2) **Magnitude of phrase command** $C^{(p)}[k]$: $Q(\Theta, \Theta')$ is maximized with respect to $C^{(p)}[k]$ when

$$C^{(p)}[k] = [\mathbf{F}_p \bar{\mathbf{x}}_p]_k (k \in \mathcal{K}_{p_0}), \quad \mathcal{K}_{p_0} = \{k | s_k = p_0\}. \quad (62)$$

- 3) **Magnitude of accent command** $C_n^{(a)}$: $Q(\Theta, \Theta')$ is maximized with respect to $C_n^{(a)}$ when

$$C_n^{(a)} = \frac{1}{|\mathcal{K}_{a_n}|} \sum_{k \in \mathcal{K}_{a_n}} [\mathbf{F}_a \bar{\mathbf{x}}_a]_k, \quad \mathcal{K}_{a_n} = \{k | s_k \in a_n\}. \quad (63)$$

- 4) **Phrase control parameter** ψ : Let us assume a Gaussian prior distribution over ψ such that

$$\psi \sim \mathcal{N}(\mu_\psi, 1/\nu_\psi^2). \quad (64)$$

Leaving only the terms in $Q(\Theta, \Theta')$ that depend on ψ , we have

$$\begin{aligned} Q(\Theta, \Theta')^\psi &= \log |\mathbf{F}_p| - \frac{1}{2} \text{tr}(\mathbf{F}_p^\top \Sigma_p^{-1} \mathbf{F}_p \mathbf{R}_p) \\ &+ \mu_p^\top \Sigma_p^{-1} \mathbf{F}_p \bar{\mathbf{x}}_p - \frac{1}{2} \nu_\psi^2 (\psi - \mu_\psi)^2. \end{aligned} \quad (65)$$

Now, let

$$\begin{aligned} U_2 &:= \begin{bmatrix} 1 & & & & O \\ -2 & 1 & & & \\ 1 & -2 & 1 & & \\ & \ddots & \ddots & \ddots & \\ O & & 1 & -2 & 1 \end{bmatrix}, U_1 := \begin{bmatrix} 0 & & & & O \\ 2 & 0 & & & \\ -2 & 2 & 0 & & \\ & \ddots & \ddots & \ddots & \\ O & & -2 & 2 & 0 \end{bmatrix}, \\ U_0 &:= \begin{bmatrix} 0 & & & & O \\ 0 & 0 & & & \\ 1 & 0 & 0 & & \\ & \ddots & \ddots & \ddots & \\ O & & 1 & 0 & 0 \end{bmatrix}, \end{aligned} \quad (66)$$

then from Eqs. (9)–(11), F_p can be written as

$$\mathbf{F}_p = \mathbf{U}_2 \psi^2 + \mathbf{U}_1 \psi + \mathbf{U}_0. \quad (67)$$

The partial derivative of $\mathcal{I}_2(\psi)$ (or $Q(\Theta, \Theta')$) with respect to ψ is a quartic function, equal up to a constant factor to

$$\begin{aligned} &2\text{tr}(\mathbf{U}_2^\top \mathbf{U}_2 \mathbf{R}_p) \psi^4 + 3\text{tr}(\mathbf{U}_2^\top \mathbf{U}_1 \mathbf{R}_p) \psi^3 \\ &+ \{\text{tr}((2\mathbf{U}_2^\top \mathbf{U}_0 + \mathbf{U}_1^\top \mathbf{U}_1) \mathbf{R}_p) - 2\mu_p^\top \mathbf{U}_2 \bar{\mathbf{x}}_p + \sigma_p^2 \nu_\psi^2\} \psi^2 \\ &+ \{\text{tr}(\mathbf{U}_1^\top \mathbf{U}_0 \mathbf{R}_p) - \mu_p^\top \mathbf{U}_1 \bar{\mathbf{x}}_p - 2\sigma_p^2 \nu_\psi^2 \mu_\psi\} \psi - 2K\sigma_p^2, \end{aligned} \quad (68)$$

and its roots, namely the stationary points of $Q(\Theta, \Theta')$, can be solved algebraically, from which we can find the optimal ψ .

- 5) **Accent control parameter** φ : Let us again assume a Gaussian prior distribution over φ such that $\varphi \sim \mathcal{N}(\mu_\varphi, 1/\nu_\varphi^2)$. As the derivation follows in exactly the same manner as above, we shall omit it.

- 6) **Baseline value** μ_b : $Q(\Theta, \Theta')$ is maximized with respect to μ_b when

$$\mu_b = \frac{\mathbf{1}^\top \Sigma_b^{-1} \bar{\mathbf{x}}_b}{\mathbf{1}^\top \Sigma_b^{-1} \mathbf{1}} = \frac{\sum_k [\bar{\mathbf{x}}_b]_k / \sigma_n[k]^2}{\sum_k 1 / \sigma_n[k]^2}. \quad (69)$$

- 7) **Variations of state emission densities** σ_p^2, σ_a^2 : $Q(\Theta, \Theta')$ is maximized with respect to $\sigma_{p,i}^2$ and $\sigma_{a,i}^2$ when

$$\sigma_p^2 = (\text{tr}(\mathbf{F}_p^\top \mathbf{F}_p \mathbf{R}_p) - 2\mu_p^\top \mathbf{F}_p \bar{\mathbf{x}}_p + \mu_p^\top \mu_p) / K, \quad (70)$$

$$\sigma_a^2 = (\text{tr}(\mathbf{F}_a^\top \mathbf{F}_a \mathbf{R}_a) - 2\mu_a^\top \mathbf{F}_a \bar{\mathbf{x}}_a + \mu_a^\top \mu_a) / K. \quad (71)$$

To summarize, we obtain the following iterative algorithm that guarantees monotonic convergence to a local maximum of the posterior density $P(\Theta | \mathbf{y})$:

- a) (E-step) Update $\bar{\mathbf{x}}_p, \bar{\mathbf{x}}_a, \mathbf{R}_p, \mathbf{R}_a$ and \mathbf{R}_b via Eqs. (54) and (55).
 - b) (M-step) Increase $Q(\Theta, \Theta')$ w.r.t. Θ through the following updates:
 - a) Update \mathbf{s} by using the Viterbi algorithm.
 - b) Update θ via Eqs. (62), (63), (70) and (71).
 - c) Update ψ (and φ) by solving the root of Eq. (68)
 - d) Update μ_b via Eq. (69).
- Return to 1) until convergence.

The complexity of this algorithm is $O(|S|^2 K)$, where $|S|$ is the number of hidden states.

B. Parameter Inference Under Non-negativity Constraints

It has been shown that phrase and accent commands must be non-negative in many non-tonal languages such as Japanese, English, German and Spanish [6]–[9]. In V-A, we treated $u_p[k]$ and $u_a[k]$ as latent variables (i.e., parameters to be marginalized out), and did not explicitly take the non-negativity constraints on $u_p[k]$ and $u_a[k]$ into consideration. While the method described in V-A can be generally applied even for such languages as Scandinavian, Portuguese and Chinese in which phrase and accent commands can be negative [10]–[13], this subsection focuses on parameter estimation under the non-negativity constraint. To impose the non-negativity constraint explicitly, it is convenient to treat $u_p[k]$ and $u_a[k]$ as model parameters instead of latent variables.

Throughout this subsection, let us assume for simplicity that ψ and φ are constants. Now, we first describe an expanded version of the generative process of \mathbf{y} :

$$\mathbf{y} | \mathbf{u}_p, \mathbf{u}_a, \mu_b \sim \mathcal{N}(\mathbf{F}_p^{-1} \mathbf{u}_p + \mathbf{F}_a^{-1} \mathbf{u}_a + \mu_b \mathbf{1}, \Sigma_b), \quad (72)$$

$$\mathbf{u}_p | \theta, \mathbf{s} \sim \mathcal{N}(\mu_p, \Sigma_p), \quad (73)$$

$$\mathbf{u}_a | \theta, \mathbf{s} \sim \mathcal{N}(\mu_a, \Sigma_a). \quad (74)$$

Note that it can be readily verified that marginalizing out \mathbf{u}_p and \mathbf{u}_a reduces Eqs. (72)–(74) to Eq. (48). Hereafter, we use \mathbf{o} to denote the set consisting of \mathbf{u}_p and \mathbf{u}_a . Instead of \mathbf{o} , we consider treating \mathbf{s} as a latent variable. Namely, we are concerned with maximizing the posterior density $P(\mathbf{o}, \theta, \mu_b | \mathbf{y}) \propto P(\mathbf{o}, \theta, \mu_b, \mathbf{y})$. We can obtain the joint probability density $P(\mathbf{o}, \theta, \mu_b, \mathbf{y})$ by marginalizing

$$P(\underbrace{\mathbf{s}, \theta, \mu_b}_{\Theta}, \mathbf{o}, \mathbf{y}) \propto \underbrace{P(\mathbf{y} | \mathbf{o}, \mu_b)}_{\text{Eq. (72)}} \underbrace{P(\mathbf{o} | \mathbf{s}, \theta)}_{\text{Eqs. (73) \& (74)}} \underbrace{P(\mathbf{s})}_{\text{Eq. (51)}}, \quad (75)$$

with respect to \mathbf{s} . We notice that the only difference from the joint probability density $P(\Theta, \mathbf{y})$ that we wanted to maximize in V-A is that \mathbf{s} is replaced with \mathbf{o} . For convenience of notation, we use $\tilde{\Theta}$ to denote the set consisting of \mathbf{o} , θ and μ_b . Hence, $\{\Theta, \mathbf{o}, \mathbf{y}\} = \{\tilde{\Theta}, \mathbf{s}, \mathbf{y}\}$.

Here we describe an iterative algorithm that searches for the maximum a posteriori estimates of $\tilde{\Theta}$ by locally maximizing $P(\tilde{\Theta}|\mathbf{y})$ given \mathbf{y} , using the generalized EM algorithm. An auxiliary function for the posterior density of interest can be written as

$$\begin{aligned} Q(\tilde{\Theta}, \tilde{\Theta}') &= \sum_{\mathbf{s} \in \mathcal{S} \times \dots \times \mathcal{S}} P(\mathbf{s}|\mathbf{y}, \tilde{\Theta}') \log P(\tilde{\Theta}, \mathbf{y}, \mathbf{s}) \\ &\stackrel{\text{def}}{=} \log P(\mathbf{y}|\mathbf{o}, \mu_b) + \sum_{\mathbf{s}} P(\mathbf{s}|\mathbf{y}, \tilde{\Theta}') \log P(\mathbf{o}|\mathbf{s}, \theta) \\ &\stackrel{\text{def}}{=} \log P(\mathbf{y}|\mathbf{o}, \mu_b) \\ &\quad + \sum_{\mathbf{s}} P(\mathbf{s}|\mathbf{y}, \tilde{\Theta}') \sum_{k=1}^K \log P(\mathbf{o}[k]|\theta, s_k) \\ &= \log P(\mathbf{y}|\mathbf{o}, \mu_b) \\ &\quad + \sum_k \sum_{s_k \in \mathcal{S}} P(s_k|\mathbf{y}, \tilde{\Theta}') \log P(\mathbf{o}[k]|\theta, s_k), \end{aligned} \quad (76)$$

where

$$\begin{aligned} \log P(\mathbf{y}|\mathbf{o}, \mu_b) &\stackrel{\text{def}}{=} - \sum_k \frac{(y[k] - x_p[k] - x_a[k] - \mu_b)^2}{2\sigma_n[k]^2}, \\ \log P(\mathbf{o}[k]|\theta, s_k) &\stackrel{\text{def}}{=} - \frac{(u_p[k] - \mu_p[k])^2}{2\sigma_p^2} - \frac{(u_a[k] - \mu_a[k])^2}{2\sigma_a^2}, \end{aligned}$$

and

$$\begin{aligned} x_p[k] &= \sum_l G_p[k-l]u_p[l], \\ x_a[k] &= \sum_l G_a[k-l]u_a[l], \\ \mu_p[k] &= [\mathbf{c}_{s_k}[k]]_1, \quad \mu_a[k] = [\mathbf{c}_{s_k}[k]]_2. \end{aligned}$$

Note that $G_p[k]$ and $G_a[k]$ stand for the discrete-time versions of $G_p(t)$ and $G_a(t)$, respectively. It can be shown that an iterative algorithm that consists of computing $\gamma = \{\gamma_{q,k}\}_{q \in \mathcal{S}, 1 \leq k \leq K}$ where $\gamma_{q,k} := P(s_k = q|\mathbf{y}, \tilde{\Theta}')$ (via the Forward-Backward algorithm), increasing $Q(\tilde{\Theta}, \tilde{\Theta}')$ with respect to $\tilde{\Theta}$, and then substituting $\tilde{\Theta}$ into $\tilde{\Theta}'$ locally maximizes the posterior $P(\tilde{\Theta}|\mathbf{y})$. Here, we must ensure that increasing $Q(\tilde{\Theta}, \tilde{\Theta}')$ with respect to \mathbf{o} is performed subject to non-negativity. This can be accomplished by invoking the idea described in [24] as follows.

By using Jensen's inequality we obtain

$$\begin{aligned} &(\hat{x}_p[k] + \hat{x}_a[k])^2 \\ &= \left(\sum_l G_p[k-l]u_p[l] + \sum_l G_a[k-l]u_a[l] \right)^2 \\ &= \left(\sum_l \lambda_{k,l}^{(p)} \frac{G_p[k-l]u_p[l]}{\lambda_{k,l}^{(p)}} + \sum_l \lambda_{k,l}^{(a)} \frac{G_a[k-l]u_a[l]}{\lambda_{k,l}^{(a)}} \right)^2 \\ &\leq \sum_l \lambda_{k,l}^{(p)} \left(\frac{G_p[k-l]u_p[l]}{\lambda_{k,l}^{(p)}} \right)^2 + \sum_l \lambda_{k,l}^{(a)} \left(\frac{G_a[k-l]u_a[l]}{\lambda_{k,l}^{(a)}} \right)^2 \\ &= \sum_l \frac{G_p[k-l]^2 u_p[l]^2}{\lambda_{k,l}^{(p)}} + \sum_l \frac{G_a[k-l]^2 u_a[l]^2}{\lambda_{k,l}^{(a)}}, \end{aligned} \quad (77)$$

where $\lambda_{k,l}^{(p)}, \lambda_{k,l}^{(a)} \geq 0$ are auxiliary variables satisfying

$$\sum_l \left(\lambda_{k,l}^{(p)} + \lambda_{k,l}^{(a)} \right) = 1. \quad (78)$$

We can verify that the equality in this inequality holds when $\lambda_{k,l}^{(p)}$ and $\lambda_{k,l}^{(a)}$ are given by

$$\hat{\lambda}_{k,l}^{(p)} = \frac{G_p[k-l]u_p[l]}{\sum_{l'} G_p[k-l']u_p[l'] + \sum_{l'} G_a[k-l']u_a[l']}, \quad (79)$$

$$\hat{\lambda}_{k,l}^{(a)} = \frac{G_a[k-l]u_a[l]}{\sum_{l'} G_p[k-l']u_p[l'] + \sum_{l'} G_a[k-l']u_a[l']}. \quad (80)$$

We can use this inequality to construct a lower bound function for $Q(\tilde{\Theta}, \tilde{\Theta}')$:

$$Q(\tilde{\Theta}, \tilde{\Theta}') \geq \check{Q}(\tilde{\Theta}, \tilde{\Theta}', \lambda), \quad (81)$$

where

$$\begin{aligned} \check{Q}(\tilde{\Theta}, \tilde{\Theta}', \lambda) &\stackrel{\text{def}}{=} - \sum_k \frac{1}{2\sigma_n[k]^2} \left\{ (y[k] - \mu_b)^2 - 2(y[k] - \mu_b)(x_p[k] + x_a[k]) \right. \\ &\quad \left. + \sum_l \left(\frac{G_p[k-l]^2 u_p[l]^2}{\lambda_{k,l}^{(p)}} + \frac{G_a[k-l]^2 u_a[l]^2}{\lambda_{k,l}^{(a)}} \right) \right\} \\ &\quad + \sum_k \sum_{s_k} p(s_k|\mathbf{y}, \tilde{\Theta}') \log p(\mathbf{o}[k]|\theta, s_k). \end{aligned} \quad (82)$$

Note that we have used λ to denote $\{\lambda_{k,l}^{(p)}, \lambda_{k,l}^{(a)}\}_{0 \leq k \leq K, 0 \leq l \leq K}$. We can use this lower bound function to derive an update equation for each element of $\tilde{\Theta}$ that guarantees a certain increase of $\check{Q}(\tilde{\Theta}, \tilde{\Theta}')$.

As mentioned above, the maximization of this lower bound function with respect to λ can be achieved when it is given by Eqs. (79) and (80). With λ fixed, the maximization of $\check{Q}(\tilde{\Theta}, \tilde{\Theta}', \lambda)$ with respect to other parameters can be achieved analytically as follows. First, $\check{Q}(\tilde{\Theta}, \tilde{\Theta}', \lambda)$ is maximized with respect to \mathbf{o} when

$$\hat{u}_p[l] = \frac{\sum_k \frac{(y[k] - \mu_b) G_p[k-l]}{\sigma_n[k]^2} + \sum_q \frac{\gamma_{q,l} [\mathbf{c}_q[l]]_1}{\sigma_p^2}}{\sum_k \frac{G_p[k-l]^2}{\sigma_n[k]^2 \lambda_{k,l}^{(p)}} + \frac{1}{\sigma_p^2}}, \quad (83)$$

$$\hat{u}_a[l] = \frac{\sum_k \frac{(y[k] - \mu_b) G_a[k-l]}{\sigma_n[k]^2} + \sum_q \frac{\gamma_{q,l} [\mathbf{c}_q[l]]_2}{\sigma_a^2}}{\sum_k \frac{G_a[k-l]^2}{\sigma_n[k]^2 \lambda_{k,l}^{(a)}} + \frac{1}{\sigma_a^2}}, \quad (84)$$

where

$$\gamma_{q,l} = p(s_l = q|\mathbf{y}, \tilde{\Theta}'). \quad (85)$$

Next, $\check{Q}(\tilde{\Theta}, \tilde{\Theta}', \lambda)$ is maximized with respect to μ_b when

$$\hat{\mu}_b = \frac{\sum_k (y[k] - x_p[k] - x_a[k]) / \sigma_n[k]^2}{\sum_k 1 / \sigma_n[k]^2}. \quad (86)$$

Finally, $\check{Q}(\tilde{\Theta}, \tilde{\Theta}', \lambda)$ is maximized with respect to $C^{(p)}[k]$, $C_n^{(a)}$, σ_p^2 and σ_a^2 when

$$\hat{C}^{(p)}[k] = u_p[k], \quad (87)$$

$$\hat{C}_n^{(p)} = \frac{\sum_{q \in a_n} \sum_k \gamma_{q,k} u_a[k] / \sigma_n[k]^2}{\sum_{q \in a_n} \sum_k \gamma_{q,k} / \sigma_n[k]^2}, \quad (88)$$

$$\hat{\sigma}_p^2 = \frac{\sum_k \sum_{q \in \mathcal{S}} \gamma_{q,k} (u_p[k] - [c_q[k]]_1)^2}{\sum_k \sum_{q \in \mathcal{S}} \gamma_{q,k}}, \quad (89)$$

$$\hat{\sigma}_a^2 = \frac{\sum_k \sum_{q \in \mathcal{S}} \gamma_{q,k} (u_a[k] - [c_q[k]]_2)^2}{\sum_k \sum_{q \in \mathcal{S}} \gamma_{q,k}}. \quad (90)$$

The above equations can be obtained by setting the partial derivatives of $\check{Q}(\tilde{\Theta}, \tilde{\Theta}', \lambda)$ with respect to $u_p[l]$, $u_a[l]$, μ_b , $C^{(p)}[k]$ and $C_n^{(a)}$ at zero, respectively. Here, it is important to note that when Eq. (83) (or Eq. (84)) becomes negative, $u_p[l] = 0$ (or $u_a[l] = 0$) is the maximizer of \check{Q} subject to non-negativity, since \check{Q} is a quadratic (strictly convex) function of $u_p[l]$ (or $u_a[l]$). Thus, under the non-negativity constraint, the update rules of $u_p[l]$ and $u_a[l]$ shall be written as

$$u_p[l] \leftarrow \max(\hat{u}_p[l], 0), \quad (91)$$

$$u_a[l] \leftarrow \max(\hat{u}_a[l], 0). \quad (92)$$

It can be shown that $Q(\tilde{\Theta}, \tilde{\Theta}')$ is non-decreasing under the updates of λ and $\tilde{\Theta}$ with the above update equations since

$$Q(\tilde{\Theta}, \tilde{\Theta}') = \check{Q}(\tilde{\Theta}, \tilde{\Theta}', \hat{\lambda}) \leq \check{Q}(\hat{\Theta}, \tilde{\Theta}', \hat{\lambda}) \leq Q(\hat{\Theta}, \tilde{\Theta}'). \quad (93)$$

To summarize, we obtain the following algorithm that guarantees the convergence to a local maximum of the posterior density of interest and the non-negativity of u_p and u_a :

- 1) Update γ using the Forward-Backward algorithm.
- 2) Increase $Q(\tilde{\Theta}, \tilde{\Theta}')$ w.r.t. $\tilde{\Theta}$ through the following updates:
 - a) Update λ via Eqs. (79) and (80).
 - b) Update \mathbf{o} via Eqs. (91) and (92).
 - c) Update θ via Eqs. (87)–(90).
 - d) Update μ_b via Eq. (86).

Return to 1) until convergence.

As with the algorithm in V-A, the complexity of this algorithm is $O(|\mathcal{S}|^2 K)$. After convergence, we search for the optimal state sequence \mathbf{s} from the output sequence $\mathbf{o} = \{\mathbf{o}[k]\}_{k=1}^K$ by using the Viterbi algorithm.

VI. EVALUATION OF PARAMETER ESTIMATION ACCURACY

A. Parameter Estimation Using Real Speech Data

To evaluate the parameter estimation accuracy of the algorithms proposed in V-A and V-B, we conducted an experiment using real speech data, excerpted from the ATR Japanese speech database B-set [25]. This database consists of 503 phonetically balanced sentences. We selected speech samples of one male speaker (MHT). The ground truth data of the Fujisaki model parameters had been manually annotated by an expert in the speech prosody field. In these ground truth data, the baseline values were all set at log 60 Hz. We chose the Fujisaki model

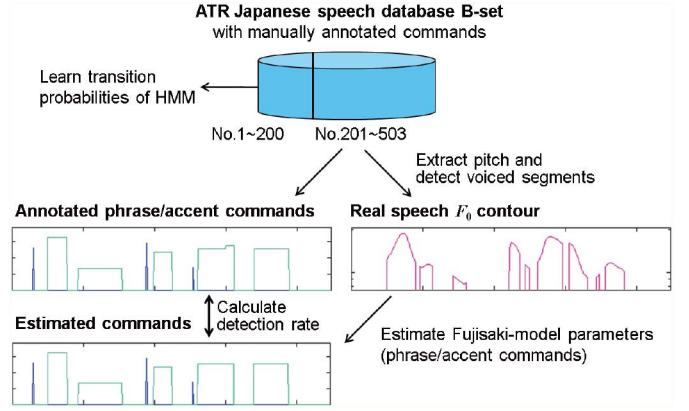


Fig. 4. Overview of the experiment in VI-A.

parameter extractor developed by Narusawa [15] as a baseline method for comparison.

Fig. 4 shows the experimental scheme of the evaluation. F_0 contours were extracted using a method we had previously developed [26], from which the Fujisaki model parameters were estimated using the present algorithm. V/UV segments were obtained by simple energy thresholding. The constant parameters were fixed at $N = 10$, $t_0 = 8$ ms, $\alpha = 3.0$ rad/s, $\beta = 20.0$ rad/s, $v_p^2 = 0.2^2$, $v_a^2 = 0.1^2$, $v_n^2[k] = 10^{15}$ for unvoiced regions and $v_n^2[k] = 0.2^2$ for voiced regions. μ_b was set at the minimum log F_0 value in the voiced regions. The initial values of Θ were set at the values obtained with Narusawa's method [15]. The EM algorithm was then run for 20 iterations. The number of substates in the HMM and the transition probability $\phi_{i',i}$ were determined according to the manually annotated data of the first 200 sentences. The parameter estimation algorithm was then tested on the remaining 303 sentences.

We evaluated the accuracy of the parameter estimation based on the following two criteria: (1) the detection rate of the phrase and accent commands, and (2) the root mean squared error (RMSE) between an observed log F_0 contour, $y[k]$, and an estimated model, $G_p[k] * \mu_p[k] + G_a[k] * \mu_a[k] + \mu_b$ over the voiced regions. The aim of this experiment was to confirm whether the present method is able to achieve accurate model fitting while ensuring that the estimated parameters are linguistically reasonable. The log F_0 RMSE indicates how well the estimated Fujisaki model fits an observed F_0 contour. The detection rate of the phrase and accent commands indicates how linguistically reasonable the estimated parameters are. The detection rate of the phrase and accent commands was calculated by matching the estimated and ground truth command sequences on a command-by-command basis using a dynamic programming algorithm. If the time difference between the estimated and ground truth phrase commands was shorter than S seconds, the estimated phrase command was considered "matched" and the local distance was set at zero. Otherwise the local distance was set at 1. As for the accent commands, we took the average of the time difference between the onsets of the estimated and ground truth accent commands and the time difference between their offsets. In the same way, when the average time difference was shorter than S seconds, the estimated accent command was considered matched. The magnitudes of the phrase and accent commands were not taken into account in our evaluation. This is because the magnitude estimation was very sensitive to the

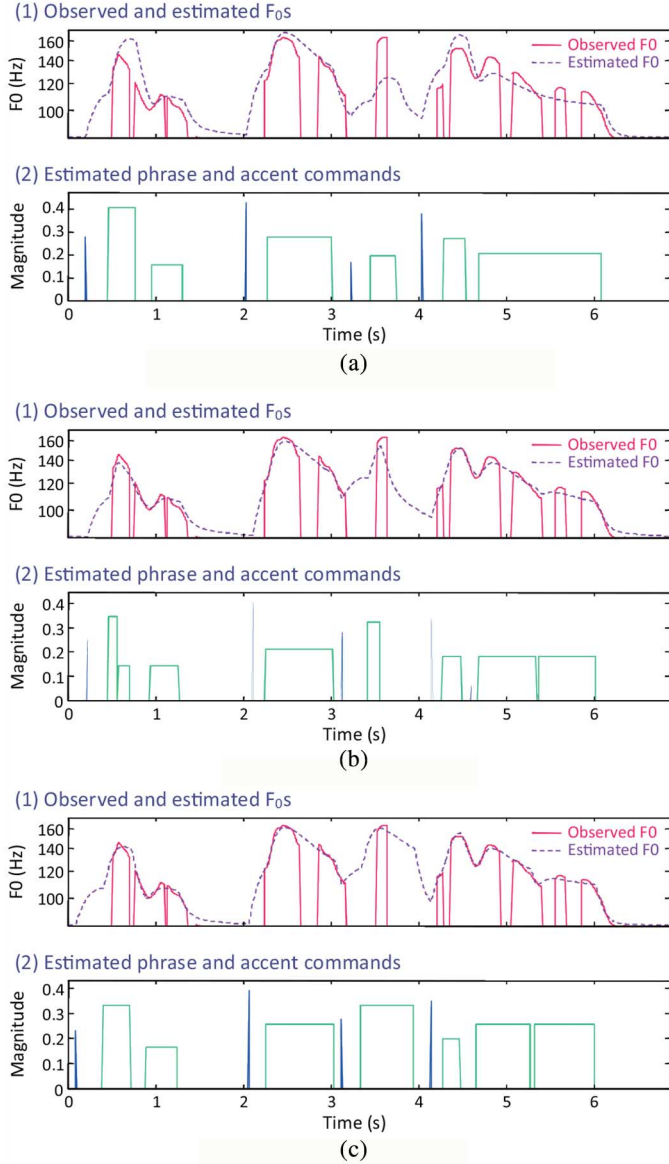


Fig. 5. Example of command detection (1 of 4). (1) An observed F_0 contour in voiced regions (solid line) and the estimated F_0 contours (dotted line) along with (2) the estimated phrase and accent commands. (a) Narusawa's method [15]. (b) Algorithm proposed in V-A. (c) Algorithm proposed in V-B.

baseline F_0 value, which was set differently in the present method and in the manual annotation. Let N_A be the number of commands in the ground truth command sequences and N_{Asum} be the sum of N_A for all 303 sentences. We defined the detection rate as

$$\frac{N_{Asum} - N_{Isum} - N_{Ssum} - N_{Dsum}}{N_{Asum}} \times 100(\%), \quad (94)$$

where N_{Isum} , N_{Ssum} and N_{Dsum} are the total numbers of insertion, substitution and deletion errors, respectively.

Table I shows the detection rate results for phrase and accent commands with $S = 0.3, 0.2, 0.1[s]$. ‘‘C,’’ ‘‘P1’’ and ‘‘P2’’ refer to Narusawa's method [15], and the algorithms proposed in V-A and V-B, respectively. The results showed that the proposed algorithms ‘‘P1’’ and ‘‘P2’’ were superior to Narusawa's method, and ‘‘P2’’ was slightly superior to ‘‘P1’’ in terms of detection rate. The left graph of Fig. 9 shows the log F_0 RMSEs. As the

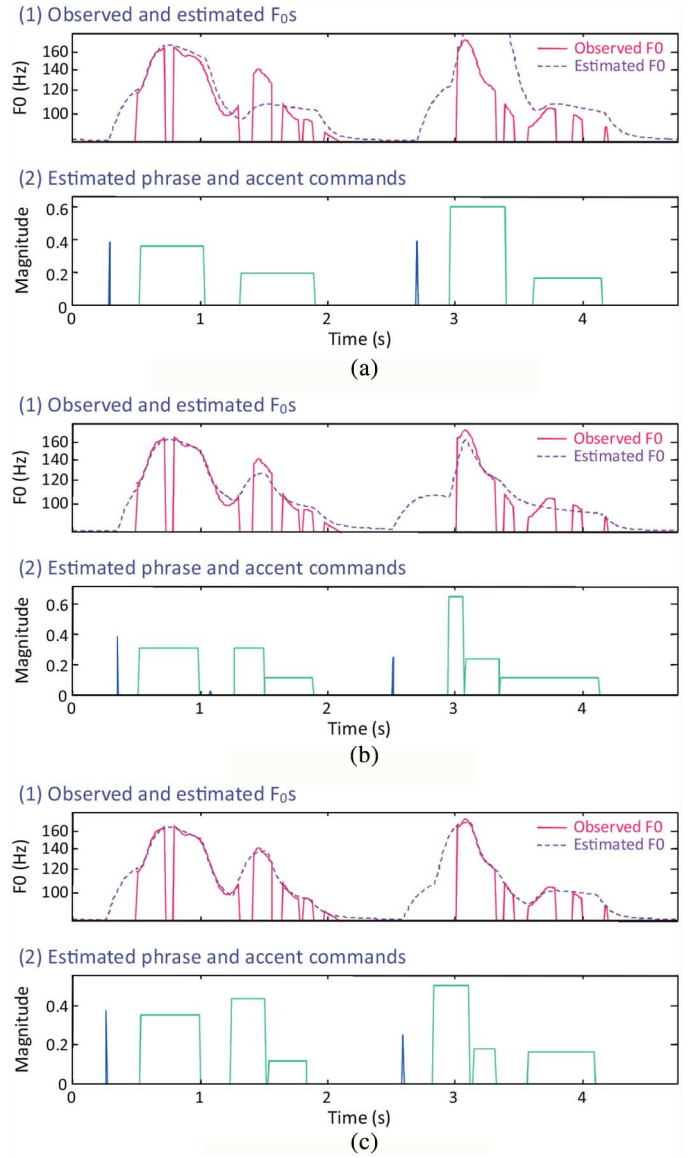


Fig. 6. Example of command detection (2 of 4). (1) An observed F_0 contour in voiced regions (solid line) and the estimated F_0 contours (dotted line) along with (2) the estimated phrase and accent commands. (a) Narusawa's method [15]. (b) Algorithm proposed in V-A. (c) Algorithm proposed in V-B.

results show, ‘‘P2’’ yielded the highest model fitting accuracy. Figs. 5–8 show some examples of observed F_0 contours and the estimated F_0 contours obtained with the conventional and present methods, from which we can confirm that the present methods (especially ‘‘P2’’) were able to fit the Fujisaki model to observed F_0 contours fairly well.

It should be noted that the detection rate tended to drop significantly when $S = 0.1[s]$. Considering the fact that average syllable durations are typically about 0.2 [s], deviations longer than 0.1 [s] from the true positions are not negligible. This implies that the proposed methods still have plenty of room for improvement.

B. Parameter Estimation Using Synthetic F_0 Contours

To evaluate the pure behavior of the present parameter estimation algorithms, we also conducted a command estimation experiment using synthetic F_0 contours. The synthetic F_0

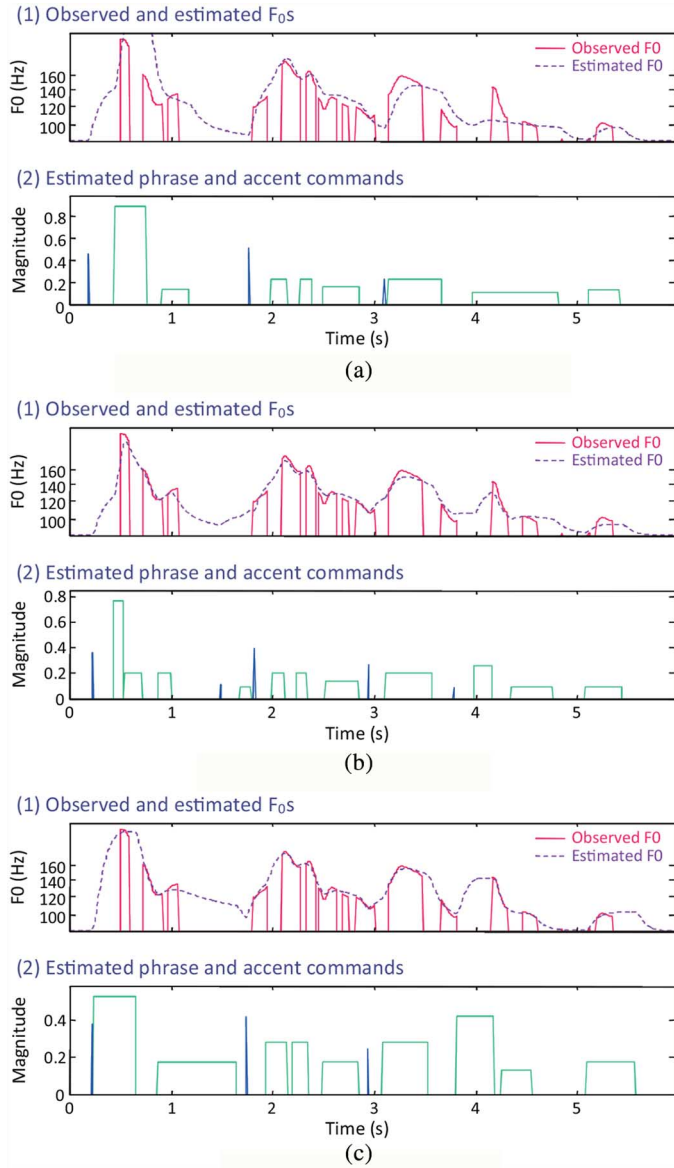


Fig. 7. Example of command detection (3 of 4). (1) An observed F_0 contour in voiced regions (solid line) and the estimated F_0 contours (dotted line) along with (2) the estimated phrase and accent commands. (a) Narusawa's method [15]. (b) Algorithm proposed in V-A. (c) Algorithm proposed in V-B.

contours were artificially created using the original Fujisaki model with the abovementioned, manually annotated command sequences. All other experimental conditions were the same as above. Fig. 10 provides an overview of this experiment.

Table II shows detection rate results for command sequences with different S settings. The $\log F_0$ RMSEs are shown in the right graph of Fig. 9. As the results show, the proposed algorithms were again significantly superior to the conventional method in terms of both the detection rate of the command sequences and the model fitting accuracy. However, as regards the accent command detection rate, the proposed methods were outperformed by the conventional method. The conventional method uses the fact that the maxima and minima of the first derivative of the F_0 contour of the Fujisaki model correspond to the onsets and offsets of accent commands with a constant delay of $1/\beta$ if the contributions of the phrase components can be disregarded. One reason why the conventional method was

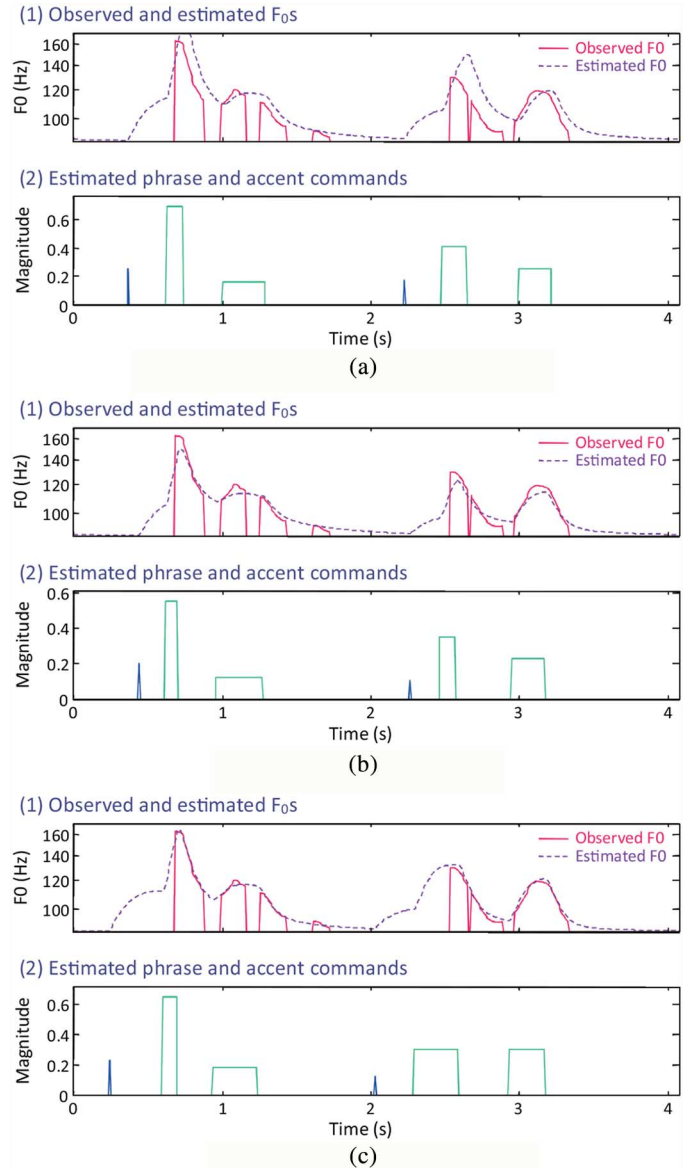


Fig. 8. Example of command detection (4 of 4). (1) An observed F_0 contour in voiced regions (solid line) and the estimated F_0 contours (dotted line) along with (2) the estimated phrase and accent commands. (a) Narusawa's method [15]. (b) Algorithm proposed in V-A. (c) Algorithm proposed in V-B.

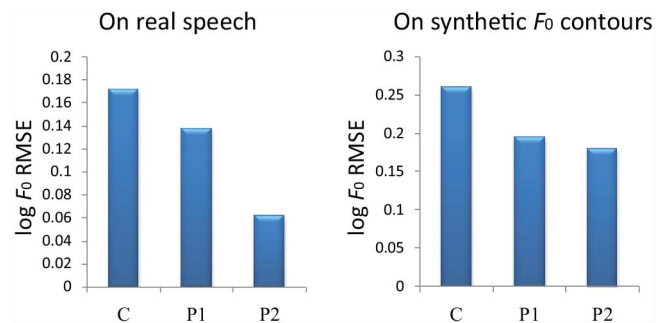


Fig. 9. The root mean squared errors (RMSEs) between observed $\log F_0$ contours and estimated models.

able to accurately detect accent commands might be that this experiment used synthetic F_0 contours created using the Fujisaki model as the test data, which agrees particularly well with the above assumption.

TABLE I
DETECTION RATE (%) OF PHRASE AND ACCENT COMMANDS WITH DIFFERENT S SETTINGS EVALUATED ON REAL SPEECH DATA

	$S = 0.3[s]$			$S = 0.2[s]$			$S = 0.1[s]$		
	C	P1	P2	C	P1	P2	C	P1	P2
All commands	89.1	91.2	93.2	84.2	85.0	88.0	64.5	58.3	67.9
Phrase commands	87.4	99.0	90.3	84.9	94.9	88.1	69.0	83.6	72.4
Accent commands	90.0	87.1	94.7	83.8	79.8	88.0	62.7	54.8	65.6

TABLE II
DETECTION RATE (%) OF PHRASE AND ACCENT COMMANDS WITH DIFFERENT S SETTINGS EVALUATED ON SYNTHETIC F_0 CONTOURS

	$S = 0.3[s]$			$S = 0.2[s]$			$S = 0.1[s]$		
	C	P1	P2	C	P1	P2	C	P1	P2
All commands	88.6	94.4	93.8	85.6	93.2	91.7	75.7	88.1	84.2
Phrase commands	71.5	99.3	95.5	66.2	98.4	91.0	52.1	97.0	85.9
Accent commands	99.3	91.8	92.8	97.5	90.5	92.1	89.9	83.5	83.4

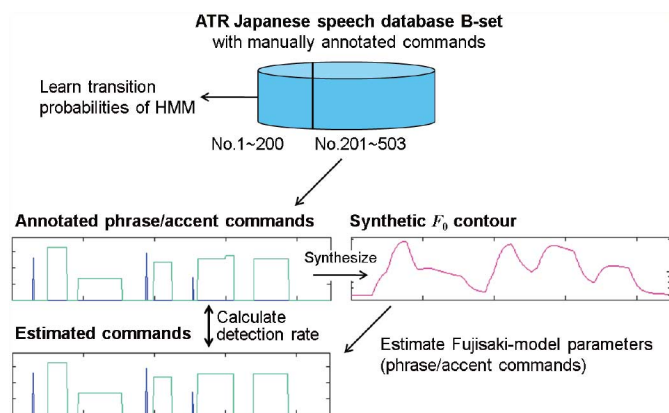


Fig. 10. Overview of the experiment described in VI-B.

VII. CONCLUSION

This paper proposed introducing a generative model of voice F_0 contours for estimating prosodic features from raw speech data. We formulated the present F_0 contour model by translating the Fujisaki model, a well-founded mathematical model representing the control mechanism of vocal fold vibration, into a probabilistic model described as a discrete-time stochastic process. There were two motivations behind this formulation. One was to derive a general parameter estimation framework for the Fujisaki model that allows the introduction of powerful algorithms such as the Viterbi algorithm, forward-backward algorithm and EM algorithm. The other was to construct an automatically trainable version of the Fujisaki model that we can naturally incorporate into statistical speech synthesis and conversion frameworks. We quantitatively evaluated the performance of the proposed Fujisaki model parameter extractor using real speech data. Experimental results revealed that our method was superior to a state-of-the-art Fujisaki model parameter extractor. The application of the present F_0 contour model to prosody generation for text-to-speech synthesis is one of our ongoing projects. A preliminary study is presented in [27].

ACKNOWLEDGMENT

We thank Prof. Keikichi Hirose (The University of Tokyo) for kindly providing us with the manually annotated data associ-

ated with the ATR speech samples. We also thank Prof. Shigeki Sagayama (Meiji University) and Dr. Daisuke Saito (The University of Tokyo) for fruitful discussions.

REFERENCES

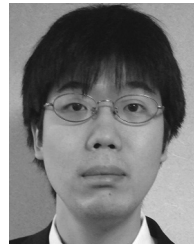
- [1] H. Kameoka, J. Le Roux, and Y. Ohishi, "A statistical model of speech F_0 contours," in *Proc. SCA Workshop Statist. Percept. Audition (SAPA'10)*, 2010, pp. 43–48.
- [2] K. Yoshizato, H. Kameoka, D. Saito, and S. Sagayama, "Statistical approach to Fujisaki-model parameter estimation from speech signals and its quantitative evaluation," in *Proc. Speech Prosody*, 2012, pp. 175–178.
- [3] K. Yoshizato, H. Kameoka, D. Saito, and S. Sagayama, "Hidden Markov convolutive mixture model for pitch contour analysis of speech," in *Proc. 13th Annu. Conf. Int. Speech Commun. Association (Interspeech'12)*, 2012.
- [4] H. Kameoka, K. Yoshizato, T. Ishihara, Y. Ohishi, K. Kashino, and S. Sagayama, "Generative modeling of speech F_0 contours," in *Proc. 14th Annu. Conf. Int. Speech Commun. Assoc. (Interspeech'13)*, 2013, pp. 1826–1830.
- [5] H. Fujisaki, O. Fujimura, Ed., "A note on the physiological and physical basis for the phrase and accent components in the voice fundamental frequency contour," in *Vocal Physiology: Voice Production, Mechanisms and Functions*. New York, NY, USA: Raven, 1988.
- [6] H. Fujisaki and K. Hirose, "Analysis of voice fundamental frequency contours for declarative sentences of Japanese," *J. Acoust. Soc. Jpn. (E)*, vol. 5, no. 4, pp. 233–242, 1984.
- [7] H. Fujisaki and S. Ohno, "Analysis and modeling of fundamental frequency contours of English utterances," in *Proc. 4th Eur. Conf. Speech Commun. Technol. (EUROSPEECH'95)*, 1995, vol. 2, pp. 985–988.
- [8] H. Mixdorff and H. Fujisaki, "Analysis of voice fundamental frequency contours of German utterances using a quantitative model," in *Proc. 3rd Int. Conf. Spoken Lang. Process. (ICSLP'94)*, 1994, vol. 4, pp. 2231–2234.
- [9] H. Fujisaki, S. Ohno, K. Nakamura, M. Guirao, and J. Gurlekian, "Analysis of accent and intonation in Spanish based on the command-response model," in *Proc. 3rd Int. Conf. Spoken Lang. Process. (ICSLP'94)*, 1994, vol. 1, pp. 355–358.
- [10] H. Fujisaki, M. Ljungqvist, and H. Murata, "Analysis and modeling of word accent and sentence intonation in Swedish," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP'93)*, 1993, pp. 211–214.
- [11] H. Fujisaki, S. Narusawa, S. Ohno, and D. Freitas, "Analysis and modeling of F_0 contours of Portuguese utterances based on the command-response model," in *Proc. 8th Eur. Conf. Speech Commun. Technol. (EUROSPEECH'03)*, 2003, vol. 3, pp. 2317–2320.
- [12] C. Wang, H. Fujisaki, R. Tomana, and S. Ohno, "Analysis of fundamental frequency contours of standard Chinese in terms of a command-response model and its application to synthesis by rule of intonation," in *Proc. 6th Int. Conf. Spoken Lang. Process. (ICSLP'00)*, 2000, vol. 3, pp. 326–329.

- [13] H. Fujisaki, W. Gu, and K. Hirose, "The command-response model for the generation of F_0 contours of Cantonese utterances," in *Proc. 7th Int. Conf. Signal Process. (ICSP'04)*, 2004, vol. 1, pp. 655–658.
- [14] H. Mixdorff, "A novel approach to the fully automatic extraction of Fujisaki model parameters," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP'00)*, 2000, vol. 3, pp. 1281–1284.
- [15] S. Narusawa, N. Minematsu, K. Hirose, and H. Fujisaki, "A method for automatic extraction of model parameters from fundamental frequency contours of speech," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP'02)*, 2002, pp. 509–512.
- [16] P. S. Rossi, F. Palmieri, and F. Cutugno, "A method for automatic extraction of Fujisaki-model parameters," in *Proc. Speech Prosody'02*, 2002, pp. 615–618.
- [17] P. S. Rossi, F. Palmieri, and F. Cutugno, "Inversion of F_0 model for natural-sounding speech synthesis," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP'03)*, 2003, pp. 520–523.
- [18] H. R. Pfizinger, H. Mixdorff, and J. Schwarz, "Comparison of Fujisaki-model extractors and F_0 stylizers," in *Proc. 10th Annu. Conf. Int. Speech Commun. Assoc. (Interspeech'09)*, 2009, pp. 2455–2458.
- [19] J. D. Ferguson, "Variable duration models for speech," in *Proc. Symp. Applicat. Hidden Markov Models to Text and Speech*, 1980, pp. 143–179.
- [20] M. Russell and R. Moore, "Explicit modeling of state occupancy in hidden Markov models for automatic speech recognition," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP'85)*, 1985, pp. 5–8.
- [21] S. Levinson, "Continuously variable duration hidden Markov models for speech analysis," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP'86)*, 1986, pp. 1241–1244.
- [22] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. R. Statist. Soc. Ser. B*, vol. 39, no. 1, pp. 1–38, 1977.
- [23] M. Feder and E. Weinstein, "Parameter estimation of superimposed signals using the EM algorithm," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 36, no. 4, pp. 477–489, Apr. 1988.
- [24] H. Kameoka, T. Nakatani, and T. Yoshioka, "Robust speech dereverberation based on non-negativity and sparse nature of speech spectrograms," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP'09)*, 2009, pp. 45–48.
- [25] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, "ATR Japanese speech database as a tool of speech recognition and synthesis," *Speech Commun.*, vol. 27, pp. 187–207, 1999.
- [26] H. Kameoka, "Statistical approach to multipitch analysis," Ph.D. dissertation, The Univ. of Tokyo, Tokyo, Japan, 2007.
- [27] K. Kadowaki, T. Ishihara, N. Hojo, and H. Kameoka, "Speech prosody generation for text-to-speech synthesis based on generative model of F_0 contours," in *Proc. 15th Annu. Conf. Int. Speech Commun. Assoc. (Interspeech'14)*, 2014, pp. 2322–2326.



Hirokazu Kameoka received B.E., M.S. and Ph.D. degrees all from the University of Tokyo, Japan, in 2002, 2004 and 2007, respectively. He is currently a Distinguished Researcher at NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation, and is also an Adjunct Associate Professor at the University of Tokyo. His research interests include audio, speech, and music signal processing and machine learning. He has served as a publicity chair of ISMIR 2009 and a program committee member of ISMIR 2014 and

ISMIR 2015. Since 2015, he has been an associate editor of the IEEE/ACM TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING. He is a member of IEEE, the Institute of Electronics, Information and Communication Engineers (IEICE), the Information Processing Society of Japan (IPJS) and the Acoustical Society of Japan (ASJ). He received 13 awards over the past 10 years, including the IEEE Signal Processing Society 2008 SPS Young Author Best Paper Award. He is the author or co-author of about 90 articles in journal papers and peer-reviewed conference proceedings.



Kota Yoshizato received the B.E. and M.S. degrees from the University of Tokyo, Japan, in 2011 and 2013, respectively. Since then, he has been a Software Engineer at DeNA Corp.



Tatsuma Ishihara received the B.E. and M.S. degrees from the University of Tokyo, Japan, in 2012 and 2014, respectively. Since then, he has been a Research Engineer with the Toshiba Corporate Research & Development Center.



Kento Kadowaki received the B.E. degree from Osaka University, Japan, in 2013 and the M.S. degree from the University of Tokyo, Japan, in 2015. He is currently with Nomura Research Institute Co., Ltd., Japan.



Yasunori Ohishi received his Ph.D. degree in information science in 2009 from Nagoya University, Japan. In 2009, he joined NTT Communication Science Laboratories. His research interests generally concern statistical acoustic signal processing, music information retrieval, and unsupervised machine learning with audio applications. In 2009, he received the 36th Awaya Kiyoshi Science Promotion Award from Acoustical Society of Japan (ASJ).



Kunio Kashino received B.E., M.S., and Ph.D. degrees all from the University of Tokyo, Japan, in 1990, 1992, and 1995, respectively. He is currently a Senior Distinguished Researcher at Nippon Telegraph and Telephone Corporation, Executive Manager of Media Information Laboratory, NTT Communication Science Laboratories, and a Visiting Professor at National Institute of Informatics. His research interests include media information processing and audio information processing. He has been involved in the research projects on computational auditory scene analysis since 1990 and robust media search since 1996.

He is a senior member of the IEEE and the Institute of Electronics, Information and Communication Engineers (IEICE), and a member of the Information Processing Society of Japan (IPJS), the Acoustical Society of Japan (ASJ), and the Japanese Society for Artificial Intelligence (JSAI). He received more than 10 awards over the past 20 years, including the IEEE TRANSACTIONS ON MULTIMEDIA Prize Paper Award.