



# Generative modeling of speech $F_0$ contours

Hirokazu Kameoka<sup>1,2</sup>, Kota Yoshizato<sup>1</sup>, Tatsuma Ishihara<sup>1</sup>,  
Yasunori Ohishi<sup>2</sup>, Kunio Kashino<sup>2</sup>, Shigeki Sagayama<sup>1</sup>

<sup>1</sup>Graduate School of Information Science and Technology, The University of Tokyo

<sup>2</sup>NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation

{kameoka, yoshizato, ishihara, sagayama}@hil.t.u-tokyo.ac.jp,  
{kameoka.hirokazu, ohishi.yasunori, kashino.kunio}@lab.ntt.co.jp

## Abstract

This paper introduces our ongoing work on generative modeling of speech fundamental frequency ( $F_0$ ) contours for estimating 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. The motivation behind this formulation is two fold. One is to derive a general parameter estimation framework for the Fujisaki model, allowing for the introduction of powerful statistical methods. The other is to construct an automatically trainable version of the Fujisaki model so that in future it can be used to develop a statistical speaking style conversion system or incorporated into existing text-to-speech synthesis systems to improve the naturalness and intelligibility of computer-generated speech. We also briefly introduce a generative model of  $F_0$  contours of singing voice developed under the same spirit.

**Index Terms:** speech  $F_0$  contour, Fujisaki model, generative model, hidden Markov model, EM algorithm

## 1. Introduction

Prosody aids the listener in interpreting 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 can be 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.

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 in phrasing while the latter contribute in 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 [1] is a well-founded mathematical model that describes an  $F_0$  contour as the sum of these two contributions. This model is known to approximate actual  $F_0$  contours of speech surprisingly well when the model parameters are chosen appropriately, and its validity has been shown for many, typologically diverse languages. For this reason, and thanks to the intuitive association of the model parameters with

the mechanical factors in the control mechanism of phonation, the Fujisaki model has been widely used with notable success to design  $F_0$  contours for synthesizing natural speech. Since a prosodic feature in speech is predominantly characterized by the levels and timings of the phrase and accent components, one important challenge is to solve an inverse problem of estimating the Fujisaki-model parameters automatically from a raw  $F_0$  contour.

However, this problem has been a difficult task. Several techniques have already been developed ([2–4], to name just a few), but so far with limited success due to the difficulty in searching for optimal parameters under the constraints imposed in the Fujisaki model. To overcome this difficulty, we have been concerned with translating the Fujisaki model into a probabilistic model such that one can make the best use of powerful methods in statistical estimation theory (such as expectation-maximization algorithm and Viterbi algorithm) for the parameter estimation. The other important motivation for this formulation is to construct an automatically trainable version of the Fujisaki model so that in future it can be used to develop a statistical speaking style conversion system or incorporated into existing text-to-speech synthesis systems to improve the naturalness and intelligibility of computer-generated speech. This paper introduces our ongoing work along with some new ideas on generative modeling of speech fundamental frequency ( $F_0$ ) contours based on a probabilistic reformulation of the Fujisaki model [5–8]. We also briefly introduce our recent work on generative modeling of  $F_0$  contours of singing voice [9] developed under the same spirit.

## 2. Original Fujisaki Model

The Fujisaki model [1], 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 component  $x_b$ :

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

The phrase component  $x_p(t)$  consists of the major-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

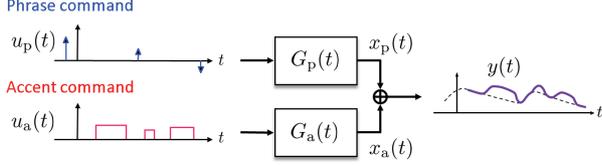


Figure 1: Original Fujisaki model [1].

of rectangular pulses (accent commands):

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

$$G_p(t) = \begin{cases} \alpha^2 t e^{-\alpha t} & (t \geq 0) \\ 0 & (t < 0) \end{cases}, \quad (3)$$

$$x_a(t) = G_a(t) * u_a(t), \quad (4)$$

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

where  $*$  denotes convolution over time. The baseline component  $x_b$  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.  $\alpha$  and  $\beta$  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.

### 3. 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 [5]. The reason for the discretization will be made apparent later. The transfer function (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  $s \simeq (1 - z^{-1})/t_0$ , 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

$$\mathcal{G}_p^{-1}(z) = a_2 z^{-2} + a_1 z^{-1} + a_0, \quad (7)$$

with  $a_2 = (\psi - 1)^2$ ,  $a_1 = -2\psi(\psi - 1)$ , and  $a_0 = \psi^2$ , where  $\psi = 1 + 1/(\alpha t_0)$ . 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. Then,  $x_p[k]$  can thus be regarded as the output of a constrained all-pole system whose characteristics are governed by a single parameter  $\psi$  (or  $\alpha$ ), such that

$$u_p[k] = a_0 x_p[k] + a_1 x_p[k - 1] + a_2 x_p[k - 2]. \quad (8)$$

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

$$u_a[k] = b_0 x_a[k] + b_1 x_a[k - 1] + b_2 x_a[k - 2], \quad (9)$$

with  $b_2 = (\varphi - 1)^2$ ,  $b_1 = -2\varphi(\varphi - 1)$ , and  $b_0 = \varphi^2$ , where  $\varphi = 1 + 1/(\beta t_0)$ . Altogether, the discrete-time version of the Fujisaki model can be expressed as the superposition of the three components:  $x_p[k] + x_a[k] + x_b$ .

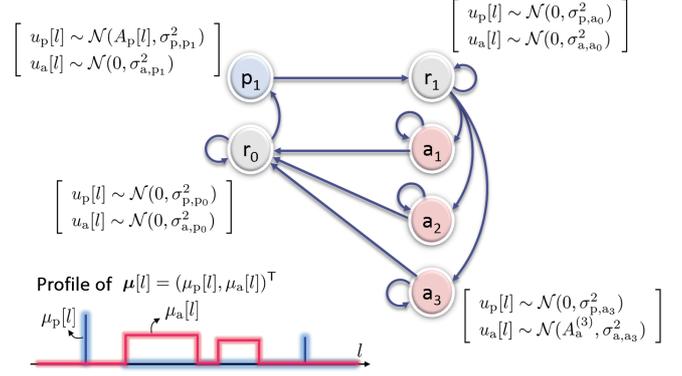


Figure 2: Command function modeling with HMM[5]. In state  $r_0$ ,  $\mu_p[k]$  and  $\mu_a[k]$  are both constrained to be zero. In state  $p_1$ ,  $\mu_p[k]$  can take a non-zero value,  $A_p[k]$ , whereas  $\mu_a[k]$  is still restricted to zero. In state  $p_1$ , no self-transitions are allowed. 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  $a_0$  leads to states  $a_1, \dots, a_N$ , in each of which  $\mu_a[k]$  can take a different non-zero value  $A_a^{(n)}$ , whereas  $\mu_p[k]$  is forced to be zero. Direct state transitions from state  $a_n$  to state  $a_{n'}$  without passing through state  $r_1$  are not allowed. This constraint restricts  $\mu_a[k]$  to consisting of rectangular pulses.

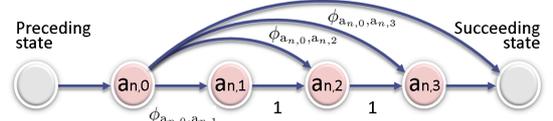


Figure 3: Duration-explicit representation of the hidden states [6]. The splitting of state  $a_n$  into substates  $a_{n,0}, a_{n,1}, a_{n,2},$  and  $a_{n,3}$  allows us to parametrize the duration of each hidden state. For example,  $\phi_{a_{n,0}, a_{n,1}}$  corresponds to the probability of staying at state  $a_n$  with 4 consecutive times.

### 4. Generative model of speech $F_0$ contours

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

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, it is required that the phrase commands must consist of Dirac deltas and the accent commands must consist of rectangular pulses. In addition, they are not allowed to overlap each other. To incorporate these requirements, we proposed in [5] to model the  $u_p[k]$  and  $u_a[k]$  pair, i.e.,  $\mathbf{o}[k] = (u_p[k], u_a[k])^T$ , using a hidden Markov model (HMM) with the specific topology illustrated in Fig. 2. The output distribution of each state is a Gaussian distribution

$$\mathbf{o}[k] \sim \mathcal{N}(\mathbf{o}[k]; \boldsymbol{\nu}[k], \boldsymbol{\Upsilon}[k]), \quad (10)$$

$$\boldsymbol{\nu}[k] = \begin{bmatrix} \mu_p[k] \\ \mu_a[k] \end{bmatrix}, \quad \boldsymbol{\Upsilon}[k] = \begin{bmatrix} \sigma_p^2[k] & 0 \\ 0 & \sigma_a^2[k] \end{bmatrix}. \quad (11)$$

where the mean vector  $\boldsymbol{\nu}[k]$  and covariance matrix  $\boldsymbol{\Upsilon}[k]$  evolve in time as a result of the state transition. To parameterize the durations of the self transitions, we proposed in [6] to split 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 flexibly control 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. Henceforth, we use the notation  $r_0 = \{r_{0,0}, r_{0,1}, \dots\}$ ,  $r_1 = \{r_{1,0}, r_{1,1}, \dots\}$ , and  $a_n = \{a_{n,0}, a_{n,1}, \dots\}$ . The HMM is defined as follows:

<p>Output sequence: <math>\{\mathbf{o}[k]\}_{k=1}^K</math>  Set of states: <math>\mathcal{S} = \{r_0, p_1, r_1, a_1, \dots, a_N\}</math>  State sequence: <math>\{s_k\}_{k=1}^K</math>  Output distribution: <math>P(\mathbf{o}[k] s_k) = \mathcal{N}(\mathbf{o}[k]; \boldsymbol{\nu}[k], \boldsymbol{\Upsilon}[k])</math>  <math display="block">\boldsymbol{\nu}[k] = \begin{cases} (0, 0)^\top &amp; (s_k \in r_0, r_1) \\ (A_p[k], 0)^\top &amp; (s_k = p_1) \\ (0, A_a^{(n)})^\top &amp; (s_k \in a_n) \end{cases}</math>  <math display="block">\boldsymbol{\Upsilon}[k] = \begin{bmatrix} \sigma_{p,s_k}^2 &amp; 0 \\ 0 &amp; \sigma_{a,s_k}^2 \end{bmatrix}</math>  Transition probability: <math>\phi_{i',i} = \log P(s_k = i   s_{k-1} = i')</math></p>
---

In [8], we further proposed designing the transition network of the HMM under the hypothesis that phrase and accent command sequences are governed by a vocabulary model. Refer to [8] for more details.

Given the state sequence  $\mathbf{s} = \{s_k\}_{k=1}^K$ , the above HMM generates the  $u_p[k]$  and  $u_a[k]$  pair. From Eq. (8) and Eq. (9),  $u_p[k]$  and  $u_a[k]$  are then fed through different all-pole systems to generate the phrase and accent components,  $x_p[k]$  and  $x_a[k]$ . It should be noted that in non-tonal languages such as standard Japanese, the phrase and accent commands must be non-negative. The treatment of the non-negativity constraint on  $x_p[k]$  and  $x_a[k]$  is discussed in [6, 7].

For real speech  $F_0$  contours, observed  $F_0$ s 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 would want to 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.,  $x_p[k] + x_a[k] + x_b$ , and a noise component  $x_n[k] \sim \mathcal{N}(0, v_n^2[k])$ , where  $v_n^2[k]$  represents the degree of uncertainty of the  $F_0$  observation at time  $k$ , which is assumed to be given. The entire  $F_0$  contour is thus given by

$$y[k] = x_p[k] + x_a[k] + x_b + x_n[k], \quad (12)$$

where  $x_b$  denotes the baseline component.

Now, let us define

$$\begin{aligned} \mathbf{u}_p &= (u_p[1], \dots, u_p[K])^\top, & \mathbf{u}_a &= (u_a[1], \dots, u_a[K])^\top, \\ \boldsymbol{\mu}_p &= (\mu_p[1], \dots, \mu_p[K])^\top, & \boldsymbol{\mu}_a &= (\mu_a[1], \dots, \mu_a[K])^\top, \\ \mathbf{x}_p &= (x_p[1], \dots, x_p[K])^\top, & \mathbf{x}_a &= (x_a[1], \dots, x_a[K])^\top, \\ \mathbf{x}_n &= (x_n[1], \dots, x_n[K])^\top, & \mathbf{y} &= (y[1], \dots, y[K])^\top. \end{aligned}$$

Then, we can write  $\mathbf{u}_p$  and  $\mathbf{u}_a$  as

$$\mathbf{u}_p = \mathbf{A}\mathbf{x}_p, \quad (13)$$

$$\mathbf{u}_a = \mathbf{B}\mathbf{x}_a, \quad (14)$$

where

$$\mathbf{A} = \begin{bmatrix} a_0 & & & & O \\ a_1 & a_0 & & & \\ a_2 & a_1 & a_0 & & \\ & \ddots & \ddots & \ddots & \\ O & & a_2 & a_1 & a_0 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} b_0 & & & & O \\ b_1 & b_0 & & & \\ b_2 & b_1 & a_0 & & \\ & \ddots & \ddots & \ddots & \\ O & & b_2 & b_1 & b_0 \end{bmatrix}. \quad (15)$$

$\mathbf{s}$  and  $\boldsymbol{\theta} = \{\{A_p[k]\}_{k=1}^K, \{A_a^{(n)}\}_{n=1}^N\}$  are the free parameters to be estimated. Obviously, estimating  $\mathbf{s}$  and  $\boldsymbol{\theta}$  corresponds to estimating the command sequences, i.e., the Fujisaki model parameters. To sum up, the likelihood function of the Fujisaki model parameters  $\mathbf{s}$  and  $\boldsymbol{\theta}$  given  $\mathbf{y}$  is given as

$$\begin{aligned} P(\mathbf{y}|\mathbf{s}, \boldsymbol{\theta}) &= \frac{|\boldsymbol{\Sigma}^{-1}|^{1/2}}{(2\pi)^{K/2}} \exp \left\{ -\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \right\}, \\ \boldsymbol{\mu} &= \mathbf{A}^{-1}\boldsymbol{\mu}_p + \mathbf{B}^{-1}\boldsymbol{\mu}_a + x_b\mathbf{1}, \\ \boldsymbol{\Sigma} &= \mathbf{A}^{-1}\boldsymbol{\Sigma}_p(\mathbf{A}^\top)^{-1} + \mathbf{B}^{-1}\boldsymbol{\Sigma}_a(\mathbf{B}^\top)^{-1} + \boldsymbol{\Sigma}_n, \end{aligned} \quad (16)$$

where

$$\boldsymbol{\Sigma}_p = \text{diag}(v_p^2[1], \dots, v_p^2[K]), \quad (17)$$

$$\boldsymbol{\Sigma}_a = \text{diag}(v_a^2[1], \dots, v_a^2[K]), \quad (18)$$

$$\boldsymbol{\Sigma}_n = \text{diag}(v_n^2[1], \dots, v_n^2[K]). \quad (19)$$

$P(\mathbf{s})$  is given by the product of the state transition probabilities:  $P(\mathbf{s}) = \phi_{s_1} \prod_{k=2}^K \phi_{s_k, s_{k-1}}$ .

Readers are referred to [5–8] for detailed derivations of the parameter inference algorithms: [5, 6] describe an iterative algorithm for maximizing the posterior density  $P(\mathbf{s}, \boldsymbol{\theta}|\mathbf{y}) \propto P(\mathbf{y}, \mathbf{s}, \boldsymbol{\theta})$  with respect to  $\mathbf{s}$  and  $\boldsymbol{\theta}$ . The devised algorithm is based on an expectation-maximization (EM) algorithm consisting in performing a Viterbi algorithm in the M-step. [7] describes a parameter estimation algorithm under non-negativity constraints on the phrase and accent components.

## 5. Singing voice $F_0$ contour modeling

So far, we have focused on the modeling of  $F_0$  contours in normal speech. Here, we briefly introduce our recent work on the generative modeling of singing voice  $F_0$  contours.

The  $F_0$  contours of a singing voice consist of two types of dynamic components. One is called a note component, which is influenced by the physical constraints of the vocal folds, such as overshoot, preparation and fine fluctuations [12, 13]. The other is called an expression component, which corresponds to a mixture of singer’s musical expressive intentions, such as vibrato and portamento [14]. Most previous papers have reported that these dynamic components strongly affect singing-voice perception, and that the former relates to the naturalness and individuality of a singing voice while the latter relates to singing styles and skills [15]. Hence, automatic decomposition of a raw  $F_0$  contour into these components can be potentially beneficial for many applications such as singer identification, singing skill evaluation, singing voice synthesis and singer style conversion.

Motivated by the above, we have been concerned with modifying the aforementioned  $F_0$  contour model so as to adapt to singing voice  $F_0$  contours. The proposed model assumes that the process of generating singing voice  $F_0$  contours is structurally similar to the Fujisaki model (Fig. 4). The note and expression components are assumed to be the outputs of second-order linear systems driven by the note and expression commands, corresponding to the musical note sequence and the musical expressive intentions, respectively. The note component contains the note transition and overshoot, whereas the expression component contains vibrato and portamento. The  $F_0$  contour is then modeled as a superposition of these two components. With the same strategy described in Sec. 4, we can translate this model into a probabilistic model and derive a parameter inference algorithm. For further details, refer to [9].

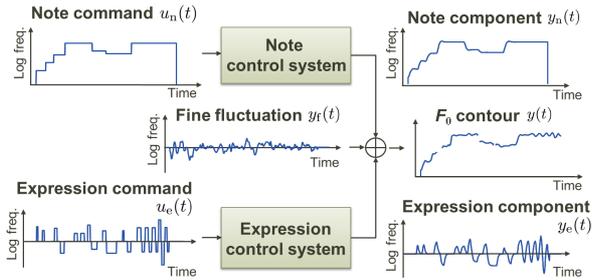


Figure 4: The proposed model of a singing voice  $F_0$  contour [9].

## 6. Experimental evaluation

We quantitatively evaluated the parameter estimation accuracy of the present algorithm using real speech data, excerpted from the ATR Japanese speech database B-set [10]. This database consists of 503 phonetically balanced sentences. We selected speech samples of one male speaker. We used Fujisaki model parameters that had been manually annotated by an expert in the field of speech prosody as the ground truth data, where the baseline component was set at  $\log 60$  Hz.  $F_0$  contours were extracted using a method we had previously developed [11], from which the Fujisaki model parameters were estimated using the present algorithm. 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:  $\log F_0$  RMSE (root mean squared error) and detection rates. Our aim was to confirm whether the present model and algorithm can achieve high model reconstruction accuracy while keeping the meaningfulness of the model parameters.  $\log F_0$  RMSE was used to evaluate the reconstruction accuracy, which measures the root mean squared error between an observed  $F_0$  contour and the estimated  $F_0$  contour. The detection rate was calculated by performing matching between the estimated and ground truth command sequences on a command-by-command basis by using a dynamic programming algorithm, where the time difference between the estimated and ground truth commands shorter than 0.3 seconds was considered “matched” and the local distance was set at zero (otherwise the local distance was set at 1). Let  $N_E$ ,  $N_A$  be the total numbers of commands in the estimated and ground truth command sequences,  $N_M$  be the number of the matched commands between the two sequences,  $N_{Esum}$ ,  $N_{Asum}$ , and  $N_{Msum}$  be the sum of  $N_E$ ,  $N_A$ ,  $N_M$  for all 303 sentences. We defined the insertion error rate  $E_I$  as  $(N_{Esum} - N_{Msum})/N_{Asum}$ , the deletion error rate  $E_D$  as  $(N_{Asum} - N_{Msum})/N_{Asum}$ , and the detection rate  $D$  as  $1 - E_I - E_D$ .

We chose Narusawa’s method [4] as a baseline method. The present method obtained a detection rate of 69.5% while the baseline method obtained 68.8%. This result confirms that our method was comparable to a state-of-the-art Fujisaki model parameter extractor in terms of the detection rate. As for the  $\log F_0$  RMSE, on the other hand, the present method obtained 0.0611, while the baseline method obtained 0.1719. This result confirms that our method was superior to the conventional method in terms of the goodness-of-fit property. Fig. 5 shows an example of the estimated phrase and accent commands obtained from a raw  $F_0$  contour.

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

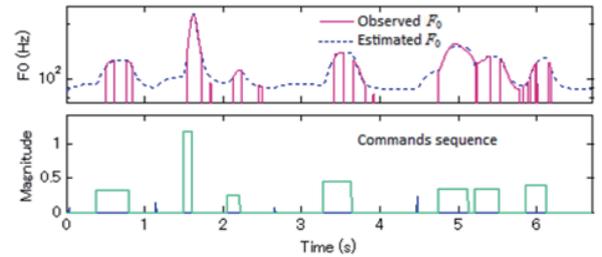


Figure 5: An example of estimated phrase and accent commands (in blue and green, respectively), along with an observed  $F_0$  contour in solid red line and the optimized  $F_0$  contour model in dotted line [7].

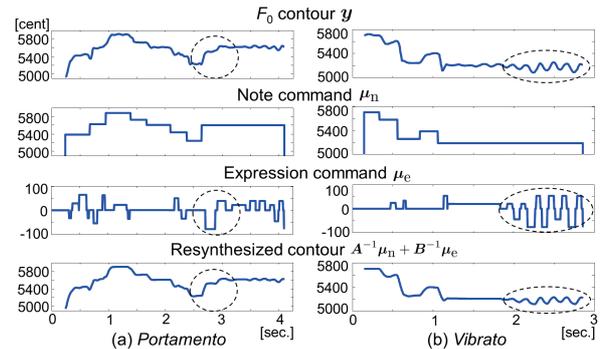


Figure 6: Estimation of note and expression commands [9].

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. The present method obtained a detection rate of 83.4% while Narusawa’s method only obtained 72.6%. With this experiment, the present method was shown to be significantly superior to the conventional method in terms of the detection rate of the command sequences.

Fig. 6 shows some examples of the decompositions of singing voice  $F_0$  contours into note and expression components by the method described in Sec. 5. As can be seen from these examples, pitch variations related to portamento and vibrato were successfully separated from an observed  $F_0$  contour.

## 7. Conclusion

This paper introduced our ongoing work on generative modeling of  $F_0$  contours in speech and singing voice based on the Fujisaki model. One important contribution of our work is that the Fujisaki model has successfully been translated into an automatically trainable model. We believe that this will open the door to applying the present model to many speech applications such as speech synthesis, speaker identification, speech conversion, and dialogue systems, in which prosodic information plays a significant role. Future work includes incorporating the present model into the HMM-based speech synthesis system (HTS) [16] in such a way that the Fujisaki-model parameters can be learned from a speech corpus in a unified manner.

## 8. Acknowledgements

We thank Prof. Keikichi Hirose (The University of Tokyo) for kindly providing us with the manually annotated data associated with the ATR speech samples, Dr. Daisuke Saito (The University of Tokyo) and Dr. Daichi Mochihashi (The Institute of Statistical Mathematics) for fruitful discussions.

## 9. References

- [1] H. Fujisaki, *In Vocal Physiology: Voice Production, Mechanisms and Functions*, Raven Press, 1988.
- [2] 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.
- [3] H. Mixdorf, "A novel approach to the fully automatic extraction of fujisaki model parameters," in *Proc. ICASSP*, 2000, vol. 3, pp. 1281–1284.
- [4] 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. ICASSP*, 2002, pp. 509–512.
- [5] H. Kameoka, J. Le Roux, and Y. Ohishi, "A statistical model of speech  $F_0$  contours," in *Proc. SAPA*, 2010, pp. 43–48.
- [6] 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*, 2012, pp. 175–178.
- [7] K. Yoshizato, H. Kameoka, D. Saito, and S. Sagayama, "Hidden Markov convolutive mixture model for pitch contour analysis of speech," in *Proc. Interspeech 2012*, 2012.
- [8] T. Ishihara, H. Kameoka, K. Yoshizato, D. Saito, and S. Sagayama, "Probabilistic speech  $F_0$  contour model incorporating statistical vocabulary model of phrase-accent command sequence," submitted to Interspeech 2013.
- [9] Y. Ohishi, H. Kameoka, D. Mochihashi, K. Kashino, "A stochastic model of singing voice  $F_0$  contours for characterizing expressive dynamic components," in *Proc. Interspeech 2012*, 2012.
- [10] 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 Communication*, vol. 27, pp. 187–207, 1999.
- [11] H. Kameoka, "Statistical speech spectrum model incorporating all-pole vocal tract model and  $F_0$  contour generating process model," in *Tech. Rep. IEICE*, 2010, in Japanese.
- [12] J. Sundberg, *The science of the singing*, Northern Illinois University Press, 1987.
- [13] G. de Krom and G. Bloothoof, "Timing and accuracy of fundamental frequency changes in singing," in *Proc. ICPhS95*.
- [14] C. E. Seashore, "A musical ornament, the vibrato," in *Psychology of Music*, 1938, pp. 33–52.
- [15] T. Saitou et al., "Acoustic and perceptual effects of vocal training in amateur male singing," in *Proc. EUROSPEECH 2009*.
- [16] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, T. Kitamura, "Speech parameter generation algorithms for HMM-based speech synthesis," in *Proc. ICASSP*, 2000, pp. 1315–1318.