

# STATISTICAL APPROACH TO FUJISAKI-MODEL PARAMETER ESTIMATION FROM SPEECH SIGNALS AND ITS QUANTITATIVE EVALUATION

Kota Yoshizato<sup>†</sup>, Hirokazu Kameoka<sup>†‡</sup>, Daisuke Saito<sup>†</sup>, Shigeki Sagayama<sup>†</sup>,

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

<sup>‡</sup> NTT Communication Science Laboratories, NTT Corporation, Japan

{yoshizato, kameoka, dsaito, sagayama}@hil.t.u-tokyo.ac.jp

## Abstract

We have previously proposed a statistical model of speech  $F_0$  contours, which is based on the discrete-time version of the Fujisaki model. One advantage of this model is that it allows us to introduce statistical methods to learn the Fujisaki-model parameters from speech  $F_0$  contours. This paper proposes several modifications to our previous model and parameter inference algorithm, and quantitatively evaluates the performance of our modified parameter inference algorithm.

**Index Terms:** Speech  $F_0$  contours, statistical model, Fujisaki model, hidden Markov model, EM algorithm

## 1. Introduction

The fundamental frequency ( $F_0$ ) contours in speech contain various types of non-linguistic information such as the speaker's identity, emotion and level of attention. They also indicate intonation in pitch accent languages. Modeling  $F_0$  contours of speech utterances is therefore potentially useful for many speech applications, in particular speech synthesis.

Thanks to the increasing availability of speech databases, speech synthesis systems based on statistical models such as hidden Markov models (HMMs) have attracted particular attention in recent years. Since unnatural  $F_0$  contours result in a synthesis that sounds "emotionless" to human listeners, one of the primary challenges in speech synthesis technology is to create a natural-sounding  $F_0$  contour for the utterance as a whole. The weakness of the current statistical text-to-speech systems is that they do not satisfactorily represent the macroscopic variations of  $F_0$ s in natural speech.

The Fujisaki model [1] is a well-founded mathematical model, which describes the process by which the whole  $F_0$  contour of a speech utterance is generated. The notable feature of the Fujisaki model is that it consists of physiologically and physically meaningful parameters (called phrase and accent commands) and is able to fit  $F_0$  contours of real speech well when they are chosen appropriately. Thus, one way of enabling statistical speech synthesizers to generate natural sounding  $F_0$  contours would be to incorporate the Fujisaki model into the statistical model so that its parameters can be learned from a speech corpus in a unified manner. However, the Fujisaki model has an analytically complex form, making it difficult to incorporate it into statistical speech synthesis systems as is.

To this end, we have recently introduced a statistical model of the discrete-time version of the Fujisaki model [2]. This model makes the best use of powerful statistical methods such as the HMM and the Expectation-Maximization (EM) algo-

rithm to estimate the Fujisaki model parameters. The aim of this paper is to (1) make several modifications to our previous model and parameter estimation algorithm, and (2) confirm the performance of our model through quantitative evaluations of its parameter estimation accuracy.

The rest of this paper is organized as follows. Section 2 briefly reviews the original Fujisaki model. Section 3 reviews our previously introduced model, i.e., a discrete counterpart of the Fujisaki model and its statistical model formulation. Section 4 presents an algorithm for estimating the Fujisaki model parameters from observed real speech  $F_0$  contour data. Section 5 shows results of a quantitative evaluation obtained by conducting an experiment using real speech data excerpted from the ATR speech database. Section 6 concludes this paper.

## 2. Original Fujisaki Model

The Fujisaki model [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  $y_p(t)$ , an accent component  $y_a(t)$ , and a base component  $y_b$ . The phrase component consists of the major-scale pitch variations over the duration of the prosodic units, and the accent component 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). The baseline component 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. The log  $F_0$  contour,  $y(t)$ , is thus expressed as

$$y(t) = y_p(t) + y_a(t) + y_b, \quad (1)$$

where

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

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

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

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

\* denotes convolution over time.  $\alpha$  and  $\beta$  are natural angular frequencies of the two second-order systems, which are known

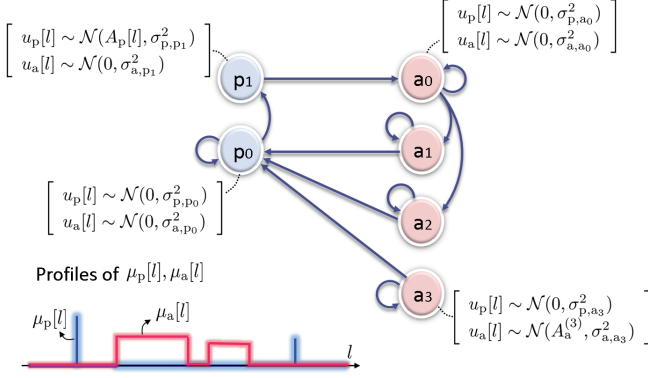


Figure 1: Command function modeling with HMM.

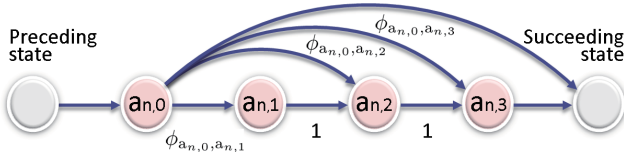


Figure 2: The splitting of state  $a_n$  into 4 substates  $a_{n,0}$ ,  $a_{n,1}$ ,  $a_{n,2}$ , and  $a_{n,3}$ .  $\phi_{a_{n,0}, a_{n,1}}$  corresponds to the probability of staying at state  $a_n$  with 4 consecutive times.

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. Statistical Speech $F_0$ Contour Model

This section describes a statistical model based on a discretized version of the Fujisaki model.

#### 3.1. Discretization

To obtain a discrete-time version of the Fujisaki model, we apply a backward difference  $s$ -to- $z$  transform to the phrase and accent control mechanisms. With this discretization, the relationships between the filter inputs and outputs are given as

$$u_p[k] = a_0 y_p[k] + a_1 y_p[k-1] + a_2 y_p[k-2], \quad (6)$$

$$a_2 = (\psi - 1)^2, \quad a_1 = -2\psi(\psi - 1), \quad a_0 = \psi^2, \quad (7)$$

$$u_a[k] = b_0 y_a[k] + b_1 y_a[k-1] + b_2 y_a[k-2], \quad (8)$$

$$b_2 = (\varphi - 1)^2, \quad b_1 = -2\varphi(\varphi - 1), \quad b_0 = \varphi^2, \quad (9)$$

where  $k$  is the discrete time index,  $y_p[k]$ ,  $u_p[k]$ ,  $y_a[k]$  and  $u_a[k]$  are the discrete-time versions of the phrase component, phrase command function, accent component and accent command function, respectively,  $\psi = 1 + 1/(\alpha t_0)$ ,  $\varphi = 1 + 1/(\beta t_0)$ , and  $t_0$  is the sampling period of the discrete-time representation.

#### 3.2. Statistical formulation

In the original Fujisaki model, the phrase commands and accent commands are assumed to consist of Dirac deltas and rectangular pulses, respectively. In addition, they are not allowed to overlap each other. To incorporate these requirements, we find it convenient 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). Specifically, we assume that  $\{\mathbf{o}[k]\}_{k=1}^K$  is a sequence of outputs

generated from an HMM with the specific topology illustrated in Figure 1. The output distribution of each state is a Gaussian distribution

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

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

where the mean vector  $\boldsymbol{\nu}[k]$  and variance matrix  $\boldsymbol{\Upsilon}[k]$  are considered to evolve in time as a result of the state transition.

To parameterize the durations of the state transitions, each state is split into a certain number of substates such that they all have exactly the same emission densities. Figure 2 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  $a_{n,l}$  to substate  $a_{n,l+1}$  is set at 1 for  $l \neq 0$ . 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  $a_{n,l}$  ( $l \geq 1$ ) corresponds to the probability of the present HMM generating a rectangular pulse that has a particular duration. In the same way, we split states  $p_0$  and  $a_0$  to parameterize the probability of the spacing between phrase and accent commands. Henceforth, we use the notation  $p_0 = \{p_{0,0}, p_{0,1}, \dots\}$ ,  $a_0 = \{a_{0,0}, a_{0,1}, \dots\}$ , and  $a_n = \{a_{n,0}, a_{n,1}, \dots\}$ . The state splitting described above is one of the modifications we have made to our previous model [2]. The present HMM is now defined as follows:

<p>Output sequence: <math>\{\mathbf{o}[k]\}_{k=1}^K</math>  Set of states: <math>\mathcal{S} = \{p_0, p_1, a_0, \dots, a_N\}</math>  State sequence: <math>\{s_k\}_{k=1}^K</math>  Output distribution: <math>P(\mathbf{o}[k]   s_k = i) = \mathcal{N}(\boldsymbol{\nu}[k], \boldsymbol{\Upsilon}[k])</math></p> $\boldsymbol{\nu}[k] = \begin{cases} (0, 0)^T & (i \in p_0, a_0) \\ (A_p[k], 0)^T & (i = p_1) \\ (0, A_a^{(n)})^T & (i \in a_n) \end{cases}$ $\boldsymbol{\Upsilon}[k] = \begin{bmatrix} \sigma_p^2 & 0 \\ 0 & \sigma_a^2 \end{bmatrix}$ <p>Transition probability: <math>\phi_{i', i} = \log P(s_k = i   s_{k-1} = i')</math></p>
---

We further assume that the baseline component is normally distributed with mean  $\mu_b$ ,  $y_b[k] \sim \mathcal{N}(\mu_b, v_b^2)$ .

For real speech  $F_0$  contours, we must take account of the uncertainty in the observed  $F_0$  data, since observed data should not always be considered reliable. For example,  $F_0$  estimates in unvoiced regions would be unreliable. Thus, we consider  $y[k]$  to be the superposition of a latent component corresponding to the true  $F_0$  and a noise component such that  $y_n[k] \sim \mathcal{N}(0, v_n^2[k])$ . The variance  $v_n^2[k]$  corresponds to the degree of uncertainty of the observed  $F_0$  at time  $k$ . Therefore,  $y[k] = y_p[k] + y_a[k] + y_b[k] + y_n[k]$ .

Now, let us define

$$\begin{aligned} \mathbf{u}_p &= (u_p[1], \dots, u_p[K])^T, & \mathbf{u}_a &= (u_a[1], \dots, u_a[K])^T, \\ \boldsymbol{\mu}_p &= (\mu_p[1], \dots, \mu_p[K])^T, & \boldsymbol{\mu}_a &= (\mu_a[1], \dots, \mu_a[K])^T, \\ \mathbf{y}_p &= (y_p[1], \dots, y_p[K])^T, & \mathbf{y}_a &= (y_a[1], \dots, y_a[K])^T, \\ \mathbf{y}_b &= (y_b[1], \dots, y_b[K])^T, & \mathbf{y}_n &= (y_n[1], \dots, y_n[K])^T, \\ \mathbf{y} &= (y[1], \dots, y[K])^T. \end{aligned}$$

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

$$\mathbf{u}_p = \mathbf{A} \mathbf{y}_p, \quad (12)$$

$$\mathbf{u}_a = \mathbf{B} \mathbf{y}_a, \quad (13)$$

where

$$\mathbf{A} = \begin{bmatrix} a_0 & & & O \\ a_1 & a_0 & & \\ a_2 & a_1 & a_0 & \\ & \ddots & \ddots & \ddots \\ O & & a_2 & a_1 & a_0 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} b_0 & & & O \\ b_1 & b_0 & & \\ b_2 & b_1 & a_0 & \\ & \ddots & \ddots & \ddots \\ O & & b_2 & b_1 & b_0 \end{bmatrix}. \quad (14)$$

For simplicity, we treat  $\phi_{i',i}$ ,  $\mu_b$ ,  $\sigma_{p,i}^2$ ,  $\sigma_{a,i}^2$ ,  $v_b^2$ ,  $v_n^2[k]$ ,  $\alpha$ ,  $\beta$  as constants and  $\Theta = \{\{A_p[k], s[k]\}_{k=1}^K, \{A_a^{(n)}\}_{n=1}^N\}$  as the free parameters to be estimated. To sum up, the likelihood function of the Fujisaki model parameters  $\Theta$  given  $\mathbf{y}$  is given as

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

$$\boldsymbol{\mu} = \mathbf{A}^{-1} \boldsymbol{\mu}_p + \mathbf{B}^{-1} \boldsymbol{\mu}_a + \mu_b \mathbf{1}, \quad (15)$$

$$\Sigma = \mathbf{A}^{-1} \Sigma_p (\mathbf{A}^\top)^{-1} + \mathbf{B}^{-1} \Sigma_a (\mathbf{B}^\top)^{-1} + \Sigma_b + \Sigma_n,$$

where

$$\Sigma_p = \begin{bmatrix} v_p^2[1] & & & O \\ & \ddots & & \\ O & & v_p^2[K] & \\ & & & \ddots \end{bmatrix}, \quad \Sigma_a = \begin{bmatrix} v_a^2[1] & & & O \\ & \ddots & & \\ O & & v_a^2[K] & \\ & & & \ddots \end{bmatrix},$$

$$\Sigma_b = \begin{bmatrix} v_b^2 & & & O \\ & \ddots & & \\ O & & v_b^2 & \\ & & & \ddots \end{bmatrix}, \quad \Sigma_n = \begin{bmatrix} v_n^2[1] & & & O \\ & \ddots & & \\ O & & v_n^2[K] & \\ & & & \ddots \end{bmatrix}. \quad (16)$$

## 4. Parameter Optimization Process

We present here an iterative algorithm that searches for the unknown parameters  $\Theta$  by locally maximizing  $P(\mathbf{y}|\Theta)$  given  $\mathbf{y}$ . By regarding  $\mathbf{x} = (\mathbf{y}_p^\top, \mathbf{y}_a^\top, \mathbf{y}_b^\top, \mathbf{y}_n^\top)^\top$  as the complete data this problem can be viewed as an incomplete data problem, which can be dealt with using the EM algorithm. The log-likelihood function of  $\Theta$  given  $\mathbf{x}$  is written as

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

$$\mathbf{x} = \begin{bmatrix} \mathbf{y}_p \\ \mathbf{y}_a \\ \mathbf{y}_b \\ \mathbf{y}_n \end{bmatrix}, \quad \mathbf{m} = \begin{bmatrix} \mathbf{A}^{-1} \boldsymbol{\mu}_p \\ \mathbf{B}^{-1} \boldsymbol{\mu}_a \\ \mu_b \mathbf{1} \\ \mathbf{0} \end{bmatrix}, \quad (17)$$

$$\Lambda^{-1} = \begin{bmatrix} \mathbf{A}^\top \Sigma_p^{-1} \mathbf{A} & \mathbf{O} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{B}^\top \Sigma_a^{-1} \mathbf{B} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \Sigma_b^{-1} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \Sigma_n^{-1} \end{bmatrix}.$$

The auxiliary function is thus given by

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

$$+ 2\mathbf{m}^\top \Lambda^{-1} \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] - \mathbf{m}^\top \Lambda^{-1} \mathbf{m} + \log P(\Theta). \quad (18)$$

$\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta]$  and  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top | \mathbf{y}; \Theta]$  are given explicitly as

$$\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta] = \mathbf{m} + \Lambda \mathbf{H}^\top (\mathbf{H} \Lambda \mathbf{H}^\top)^{-1} (\mathbf{y} - \mathbf{H} \mathbf{m}), \quad (19)$$

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

$$+ \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta] \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta]^\top, \quad (20)$$

by using the relationship  $\mathbf{y} = \mathbf{H}\mathbf{x}$ , where  $\mathbf{H} = [\mathbf{I}, \mathbf{I}, \mathbf{I}, \mathbf{I}]$ . These are the values to be updated at the E step. Let  $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta]$

be partitioned into four  $K \times 1$  blocks such that  $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta] = (\bar{\mathbf{x}}_p^\top, \bar{\mathbf{x}}_a^\top, \bar{\mathbf{x}}_b^\top, \bar{\mathbf{x}}_n^\top)^\top$ .

The M step update formulas are given as follows.

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

$$\mathcal{I}(s) := -\frac{1}{2} \sum_{k=1}^K (\mathbf{o}[k] - \boldsymbol{\nu}[k])^\top \boldsymbol{\Upsilon}[k]^{-1} (\mathbf{o}[k] - \boldsymbol{\nu}[k])$$

$$+ \log P(s_1) + \sum_{k=2}^K \log P(s_k | s_{k-1}), \quad (21)$$

where  $\mathbf{o}[k] := ([\mathbf{A}\bar{\mathbf{x}}_p]_k, [\mathbf{B}\bar{\mathbf{x}}_a]_k)^\top$ . Here the notation  $[\cdot]_k$  is used to denote the  $k$ -th element of a vector. The state sequence  $\{s_k\}_{k=1}^K$  maximizing  $\mathcal{I}(s)$  can be solved efficiently using the Viterbi algorithm.

**2) Magnitudes of phrase and accent commands:**  $Q(\Theta, \Theta')$  is maximized with respect to  $A_p[k]$  and  $A_a^{(n)}$  when

$$A_p[k] = [\mathbf{A}\bar{\mathbf{x}}_p]_k \quad (k \in \mathcal{T}_{p_1}), \quad \mathcal{T}_{p_1} = \{k | s_k = p_1\}, \quad (22)$$

$$A_a^{(n)} = \frac{1}{|\mathcal{T}_{a_n}|} \sum_{k \in \mathcal{T}_{a_n}} [\mathbf{B}\bar{\mathbf{x}}_a]_k, \quad \mathcal{T}_{a_n} = \{k | s_k = a_n\}. \quad (23)$$

The M step described above involves the computation of  $\mathbf{A}\bar{\mathbf{x}}_p$  and  $\mathbf{B}\bar{\mathbf{x}}_a$ . This amounts to computing the filter inputs  $\bar{\mathbf{u}}_p$  and  $\bar{\mathbf{u}}_a$  from the estimates of the phrase and accent components,  $\bar{\mathbf{x}}_p$  and  $\bar{\mathbf{x}}_a$ , using the ‘‘inverse filtering’’ matrices  $\mathbf{A}$  and  $\mathbf{B}$ . Note that as  $\bar{\mathbf{u}}_p$  and  $\bar{\mathbf{u}}_a$  correspond to the estimates of command functions, they should be non-negative. However, the computations of  $\mathbf{A}\bar{\mathbf{x}}_p$  and  $\mathbf{B}\bar{\mathbf{x}}_a$  allow each element of the resulting vector to have a negative value. To effectively avoid allowing such unwanted estimates, we consider finding  $\bar{\mathbf{u}}_p$  and  $\bar{\mathbf{u}}_a$  that minimize  $\|\mathbf{A}^{-1}\bar{\mathbf{u}}_p - \bar{\mathbf{x}}_p\|_2^2$  and  $\|\mathbf{B}^{-1}\bar{\mathbf{u}}_a - \bar{\mathbf{x}}_a\|_2^2$  subject to non-negativity, instead of simply computing  $\mathbf{A}\bar{\mathbf{x}}_p$  and  $\mathbf{B}\bar{\mathbf{x}}_a$ . In elementwise notation, they can be written as  $\sum_k |G_p[k] * \bar{u}_p[k] - \bar{x}_p[k]|^2$  and  $\sum_k |G_a[k] * \bar{u}_a[k] - \bar{x}_a[k]|^2$ , where  $G_p[k]$  and  $G_a[k]$  are the discrete-time versions of the impulse responses of the phrase and accent control mechanisms. Fortunately, this non-negative deconvolution problem can be solved efficiently by employing the method described in [3].

## 5. Experiment

One important contribution of our work is that the Fujisaki model has successfully been translated into a statistical model. We believe that this will open the door to combining our model and the HMM-based text-to-speech synthesis system [4] so that the Fujisaki-model parameters as well as the spectral parameter sequences can be learned from a speech corpus in a unified manner. In this regard, our model is already superior to conventional ‘‘non-statistical’’ methods such as [5]. However, it is not yet clear whether our statistical model is able to estimate the Fujisaki model parameters from real speech data as accurately as the state-of-the-art technique [5]. Thus, we quantitatively evaluated the parameter estimation accuracy of the present algorithm using real speech data, excerpted from the ATR Japanese speech database B-set [6]. This database consists of 503 phonetically balanced sentences. We selected speech samples of one male speaker (MHT). 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 constant at log 60 Hz.

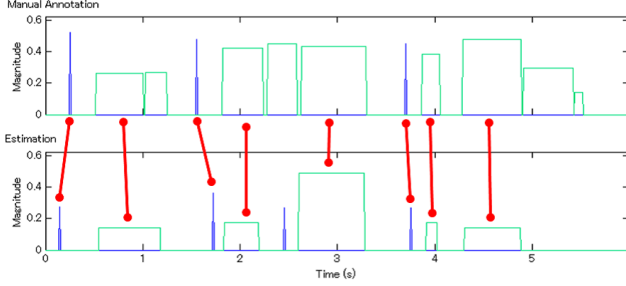


Figure 3: An example of command sequence matching.

Table 1: Accuracy rate and Error rate ( $S=0.3s$ ).

All Commands	$A$	$E_I$	$E_D$
Init	0.688	0.088	0.224
Estimated	0.697	0.127	0.177
Phrase Commands	$A$	$E_I$	$E_D$
Init	0.647	0.109	0.244
Estimated	0.680	0.207	0.112
Accent Commands	$A$	$E_I$	$E_D$
Init	0.711	0.076	0.213
Estimated	0.708	0.083	0.207

$F_0$  contours were extracted using the method described in [7], from which the Fujisaki model parameters were estimated using the present algorithm. The constant parameters were fixed respectively at  $N = 20$ ,  $t_0 = 8$  ms,  $\alpha = 3.0$  rad/s,  $\beta = 20.0$  rad/s,  $v_p^2[k] = 0.2^2$ ,  $v_a^2[k] = 0.02^2$ ,  $v_b^2 = 0.001^2$ ,  $v_n^2[k] = 10^{15}$  for unvoiced regions and  $v_n^2[k] = 10^{-15}$  for voiced regions.  $\mu_b$  was set at the minimum  $\log F_0$  value in the voiced regions. The initial values of  $\Theta$  were set at the values obtained with the method described in [5]. The EM algorithm was then run for 20 iterations. The number of substates in the present 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 parameter estimation accuracy in the following manner: We performed matching between the estimated and ground truth command sequences as illustrated in Figure 3 on a command-by-command basis using dynamic programming. 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 the offsets of the estimated and ground truth accent commands. 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 baseline  $F_0$  value, which was set differently in the present method and in the manual annotation. 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 accuracy rate  $A$  as  $1 - E_I - E_D$ . Table 1 shows the result of our quantitative evaluation with  $S = 0.3$  s. The top, middle, and bottom tables show the accuracy and error rates of the phrase and accent commands, the phrase commands alone, and the accent commands alone, respectively. The “Init” row shows the accuracy and error rates of the initial command sequence (which was obtained with the method described in [5]), and the “Estimated” row shows that of the estimated command sequence after the EM iterations. As the result shows, our method performed slightly better than the state-of-the-art technique [5].

## 6. Conclusion

In this paper, we made several modifications to our previously reported model of speech  $F_0$  contours and parameter estimation algorithm. We evaluated the parameter estimation accuracy of the present method using real speech data. Future work will include incorporating the present model into the HMM-based speech synthesis system (HTS) [4] in such a way that the Fujisaki-model parameters can be learned from a speech corpus in a unified manner. We are also currently exploring a reasonable model for the  $F_0$  contours of a singing voice [8].

## 7. Acknowledgement

We thank Prof. Keikichi Hirose of the University of Tokyo, who kindly provided us with the manually annotated ground truth data associated with the ATR speech samples.

## 8. References

- [1] H. Fujisaki, *In Vocal Physiology: Voice Production, Mechanisms and Functions*, Raven Press, 1988.
- [2] H. Kameoka, J. L. Roux, and Y. Ohishi, “A statistical model of speech  $F_0$  contours,” in *Proc. ISCA Tutorial and Research Workshop on Statistical And Perceptual Audition (SAPA)*, 2010, pp. 43–48.
- [3] H. Kameoka, T. Nakatani, and T. Yoshioka, “Robust speech dereverberation based on non-negativity and sparse nature of speech spectrograms,” in *In Proc. 2008 IEEE Intl. Conf. Acoust., Speech, Signal Process. (ICASSP 2009)*, 2009, pp. 45–48.
- [4] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, “Speech parameter generation algorithms for HMM-based speech synthesis,” in *Proc. 2000 IEEE Intl. Conf. Acoust., Speech, Signal Process. (ICASSP 2000)*, 2000, pp. 1315–1318.
- [5] 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. 2002 IEEE Intl. Conf. Acoust., Speech, Signal Process. (ICASSP 2002)*, 2002, pp. 509–512.
- [6] 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.
- [7] H. Kameoka, “Statistical speech spectrum model incorporating all-pole vocal tract model and  $F_0$  contour generating process model,” in *Technical report of the Institute of Electronics, Information and Communication Engineers (IEICE)*, 2010, in Japanese.
- [8] Y. Ohishi, H. Kameoka, D. Mochihashi, H. Nagano, and K. Kashino, “Statistical modeling of  $F_0$  dynamics in singing voices based on Gaussian processes with multiple oscillation bases,” in *Proc. 11th Annual Conference of the International Speech Communication Association (INTERSPEECH 2010)*, 2010, pp. 2598–2601.