

Hidden Markov Convulsive Mixture Model for Pitch Contour Analysis of Speech

Kota Yoshizato¹, Hirokazu Kameoka^{1,2}, Daisuke Saito¹, Shigeki Sagayama¹,

¹Graduate School of Information Science and Technology, The University of Tokyo, Japan

² NTT Communication Science Laboratories, NTT Corporation, Japan

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

Abstract

This paper proposes a stochastic model of speech F_0 contours, based on the stochastic formulation of the Fujisaki model. Our motivation for the stochastic formulation is twofold. Firstly, it allows us to derive a well-behaved algorithm for estimating the Fujisaki model parameters from a raw F_0 contour. Secondly, it will open the door to incorporating the well-founded F_0 contour model into various statistical speech processing problems. We quantitatively evaluated the performance of our method in terms of an Fujisaki-model parameter estimation accuracy using real speech data. Experimental results revealed that our method was superior to a state-of-the-art Fujisaki model parameter extractor.

Index Terms: speech F_0 contours, statistical model, Fujisaki model, hidden Markov model, EM algorithm

1. Introduction

The fundamental frequency (F_0) contours in normal speech contains various types of non-linguistic information such as speaker's identity, emotion and attention. Modeling the F_0 contours of speech utterances can thus be potentially very useful for many speech applications, including speech recognition, speaker recognition, speech synthesis, and dialogue systems.

The Fujisaki model [1] is a well-founded mathematical model, which describes the generating process of the whole F_0 contour of a speech utterance. The remarkable feature of the Fujisaki model is that it consists of physiologically and physically meaningful parameters (called the phrase and accent commands) and is able to fit F_0 contours of real speech well when they are chosen appropriately. For this reason, the Fujisaki model has often been used to manually design an F_0 contour for synthesizing natural speech. Thus, to enable speech synthesizers to automatically generate natural-sounding F_0 contours, one way would be to incorporate the Fujisaki model into the generative model of phonemic sequences so that its parameters can be learned from a speech corpus in a unified manner. However, estimating (learning) the Fujisaki model parameters from raw F_0 contour observations has been a difficult task. Several techniques have already been developed [2, 3, 4], for the purpose of incorporating the extracted parameters in automatic speech/emotion recognition systems to improve their per-

formance, but so far with limited success due to the analytical complexity of the Fujisaki model.

We have previously derived a stochastic model of speech F_0 contours by translating the Fujisaki model into a probabilistic generative model [5, 6]. The stochastic reformulation of the Fujisaki model has allowed us to derive an efficient algorithm for automatically estimating the Fujisaki-model parameters from raw F_0 contours. In this paper, we present an improved version of our previous models described in [5, 6]. The rest of this paper is organized as follows. Section 2 briefly reviews the original Fujisaki model. Section 3 formulates a probabilistic generative model of speech F_0 contours based on the Fujisaki model. Section 4 derives an algorithm for finding the maximum a posteriori (MAP) estimates of the Fujisaki model parameters from an observed F_0 contour. Section 5 presents results of a quantitative evaluation conducted 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 :

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

The phrase component $y_p(t)$ consists of the major-scale pitch variations over the duration of the prosodic units, and the accent component $y_a(t)$ 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):

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

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

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

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

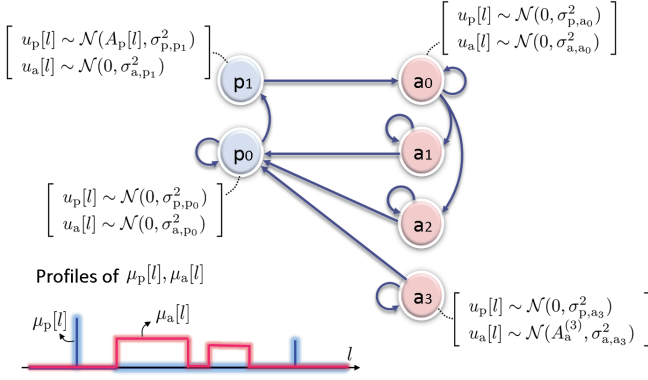


Figure 1: Command function modeling with HMM.

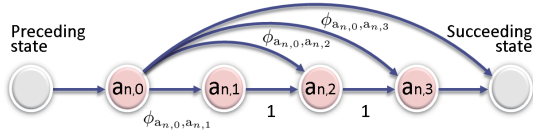


Figure 2: The splitting of state a_n into 4 substates $a_{n,0}$, $a_{n,1}$, $a_{n,2}$, and $a_{n,3}$. $\phi_{a_{n,0}, a_{n,1}}$ corresponds to the probability of staying at state a_n with 4 consecutive times.

where $*$ denotes convolution over time. The baseline component y_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.

3. Stochastic model of speech F_0 contours

Here, we model the generative process of an entire F_0 contour of speech 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 $u_a[k]$, where k denotes the discrete-time index. 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 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 emitted from an HMM with the specific topology illustrated in Fig. 1. The output distribution of each state is a Gaussian distribution

$$\mathbf{o}[k] \sim \mathcal{N}(\mathbf{o}[k]; \boldsymbol{\nu}[k], \boldsymbol{\Upsilon}[k]), \quad (6)$$

$$\boldsymbol{\nu}[k] = \begin{bmatrix} \mu_p[k] \\ \mu_a[k] \end{bmatrix}, \quad \boldsymbol{\Upsilon}[k] = \begin{bmatrix} \nu_p^2[k] & 0 \\ 0 & \nu_a^2[k] \end{bmatrix}, \quad (7)$$

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 tran-

sition. To parameterize the durations of the self transitions, each state is split into a certain number of substates such that they all have exactly the same emission densities. Fig. 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 $a_{n,m}$ 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 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 present HMM is now defined as follows:

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

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 (2) and (4), $u_p[k]$ and $u_a[k]$ are then fed through different critically damped filters, $G_p[k]$ and $G_a[k]$, to generate the phrase and accent components, $x_p[k]$ and $x_a[k]$:

$$x_p[k] = u_p[k] * G_p[k], \quad (8)$$

$$x_a[k] = u_a[k] * G_a[k], \quad (9)$$

where $*$ denotes convolution over k . The entire F_0 contour is then given by

$$x[k] = x_p[k] + x_a[k] + u_b, \quad (10)$$

where u_b denotes the baseline component. In non-tonal languages such as standard Japanese, the phrase and accent commands should be non-negative. In our previous model [5, 6], we treated $u_p[k]$, $u_a[k]$ and u_b 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. To explicitly impose the non-negativity constraints, it is convenient to treat these variables as model parameters instead of latent variables.

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 $x[k]$ 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.

Overall, an observed F_0 contour $y[k]$ is described as $y[k] = x[k] + x_n[k]$. For simplicity, we henceforth treat $\phi_{i',i}$, u_b , $\sigma_{p,i}^2$, $\sigma_{a,i}^2$, $v_n^2[k]$, α , and β as constants. By marginalizing $x_n[k]$ out, we obtain the probability density function of $\mathbf{y} = \{y[k]\}_{k=1}^K$, given $\mathbf{o} = \{o[k]\}_{k=1}^K$, as

$$P(\mathbf{y}|\mathbf{o}) = \prod_{k=1}^K \mathcal{N}(y[k]; x[k], v_n^2[k]),$$

$$x[k] = G_p[k] * u_p[k] + G_a[k] * u_a[k] + u_b. \quad (11)$$

Recall from (6) that given a state sequence $\mathbf{s} = \{s_k\}_{k=1}^K$ and $\boldsymbol{\theta} = \{\{A_p[k]\}_{k=1}^K, \{A_a^{(n)}\}_{n=1}^N\}$, \mathbf{o} is generated according to $P(\mathbf{o}|\mathbf{s}, \boldsymbol{\theta}) = \prod_{k=1}^K \mathcal{N}(o[k]; \nu[k], \Upsilon[k])$. $P(\mathbf{s})$ is given by the product of the state transition probabilities: $P(\mathbf{s}) = \phi_{s_1} \prod_{k=2}^K \phi_{s_k, s_{k-1}}$. Furthermore, we assume that $\boldsymbol{\theta}$ is uniformly distributed.

4. Parameter Optimization Process

In this section, we describe an iterative algorithm that searches for the maximum a posteriori estimates of \mathbf{o} and $\boldsymbol{\theta}$ by locally maximizing $P(\mathbf{o}, \boldsymbol{\theta}|\mathbf{y})$ given \mathbf{y} using the generalized Expectation-Maximization (EM) algorithm. We treat \mathbf{s} as a latent variable and consider marginalizing $P(\mathbf{o}, \boldsymbol{\theta}, \mathbf{s}|\mathbf{y}) \propto P(\mathbf{y}|\mathbf{o})P(\mathbf{o}|\mathbf{s}, \boldsymbol{\theta})P(\mathbf{s})$ with respect to \mathbf{s} to obtain the objective $P(\mathbf{o}, \boldsymbol{\theta}|\mathbf{y})$. The auxiliary function (as known as the “Q-function”) can be written as

$$Q(\mathbf{o}, \boldsymbol{\theta}, \mathbf{o}', \boldsymbol{\theta}') = \sum_{\mathbf{s}} P(\mathbf{s}|\mathbf{y}, \mathbf{o}', \boldsymbol{\theta}') \log P(\mathbf{o}, \boldsymbol{\theta}, \mathbf{s}|\mathbf{y})$$

$$\stackrel{c}{=} \log P(\mathbf{y}|\mathbf{o}) + \sum_{\mathbf{s}} P(\mathbf{s}|\mathbf{y}, \mathbf{o}', \boldsymbol{\theta}') \log P(\mathbf{o}|\mathbf{s}, \boldsymbol{\theta})P(\mathbf{s}),$$

where $\stackrel{c}{=}$ denotes equality up to constant terms. An iterative algorithm that consists of computing $P(\mathbf{s}|\mathbf{y}, \mathbf{o}', \boldsymbol{\theta}')$ (via the Forward-Backward algorithm), increasing $Q(\mathbf{o}, \boldsymbol{\theta}, \mathbf{o}', \boldsymbol{\theta}')$ with respect to \mathbf{o} and $\boldsymbol{\theta}$, and then substituting \mathbf{o} and $\boldsymbol{\theta}$ into \mathbf{o}' and $\boldsymbol{\theta}'$ locally maximizes the posterior $P(\mathbf{o}, \boldsymbol{\theta}|\mathbf{y})$. Here, care must be taken that increasing $Q(\mathbf{o}, \boldsymbol{\theta}, \mathbf{o}', \boldsymbol{\theta}')$ with respect to \mathbf{o} must be performed subject to non-negativity. This can be done by invoking the idea of [7]. By using the Jensen’s inequality we obtain an inequality

$$- \left(\sum_{i \in \{p,a,b\}} \sum_l G_i[k-l] u_i[l] \right)^2$$

$$\geq - \sum_{i \in \{p,a,b\}} \sum_l \frac{G_i^2[k-l] u_i^2[l]}{\lambda_{i,k,l}}, \quad (12)$$

where $G_b[k] = \delta[k]$ (Kronecker’s delta), $\lambda_{i,k,l} \geq 0$ is an auxiliary variable satisfying $\sum_i \sum_l \lambda_{i,k,l} = 1$. We can use this inequality to construct a lower bound function

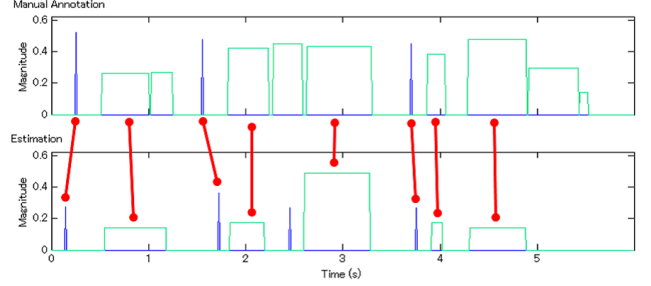


Figure 3: An example of command sequence matching.

for $Q(\mathbf{o}, \boldsymbol{\theta}, \mathbf{o}', \boldsymbol{\theta}')$. The maximization of this lower bound function w.r.t. \mathbf{o} (subject to non-negativity) and λ can be achieved analytically, which guarantees a certain increase of $Q(\mathbf{o}, \boldsymbol{\theta}, \mathbf{o}', \boldsymbol{\theta}')$.

After convergence, we search for the optimal state sequence \mathbf{s} by using the Viterbi algorithm.

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 various statistical speech applications 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 [4]. 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. 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 [8]. 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 at $\log 60$ Hz. F_0 contours were extracted using a method we had previously developed [9], from which the Fujisaki model parameters were estimated using the present algorithm. The constant parameters were fixed respectively at $N = 10$, $t_0 = 8$ ms, $\alpha = 3.0$ rad/s, $\beta = 20.0$ rad/s, $v_p^2[k] = 0.2^2$, $v_a^2[k] = 0.1^2$, $v_b^2 = 0.001^2$, $v_n^2[k] = 10^{15}$ 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 the non-statistical method [4]. 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: $\log F_0$ 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

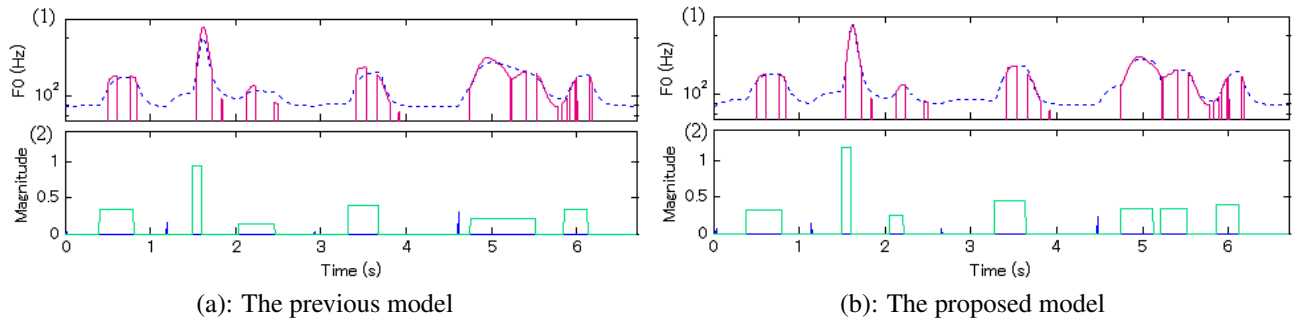


Figure 4: (1) An observed F_0 contour in voiced regions (in solid line) and the estimated F_0 contours (in dotted line) along with (2) the estimated phrase and accent commands.

Table 1: Detection rates and $\log F_0$ RMSE ($S=0.3s$).

	Detection rates	$\log F_0$ [Hz] RMSE
Init	0.688	0.1719
Estimated	0.695	0.0611

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 used to evaluate the meaningfulness of the parameter estimates, which was calculated in the following way: We performed matching between the estimated and ground truth command sequences as illustrated in Fig. 3 on a command-by-command basis by 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 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 detection rate D as $1 - E_I - E_D$.

Tab. 1 shows the result of our quantitative evaluation with $S = 0.3$ s. The “Init” row shows the detection rate and $\log F_0$ RMSE of the initial command sequence (which was obtained with the non-statistical method [4]), and the “Estimated” row shows that of the estimated com-

mand sequence after the EM iterations. From the results, we confirmed that our method was comparable to a state-of-the-art Fujisaki model extractor in terms of the detection rate. On the other hand, our method was superior to the conventional method in terms of the model reconstruction accuracy. We can also confirm from Fig. 4 that the present model is able to fit an observed F_0 contour more accurately than our previous model.

6. Conclusion

In this paper, we proposed a statistical model of speech F_0 contours and parameter estimation algorithm. We evaluated the parameter estimation accuracy of the proposed method using real speech data, and confirm the advantage of the proposed method. Future work will include incorporating the present model into the statistical speech applications such as the HMM-based speech synthesis system (HTS) in such a way that the Fujisaki-model parameters can be learned from a speech corpus in a unified manner.

7. 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. L. 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] H. Kameoka, T. Nakatani, and T. Yoshioka, “Robust speech dereverberation based on non-negativity and sparse nature of speech spectrograms,” in *Proc. ICASSP, 2009*, pp. 45–48.
- [8] 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.
- [9] 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.