

# STATISTICAL $F_0$ PREDICTION FOR ELECTROLARYNGEAL SPEECH ENHANCEMENT CONSIDERING GENERATIVE PROCESS OF $F_0$ CONTOURS WITHIN PRODUCT OF EXPERTS FRAMEWORK

Kou Tanaka<sup>1</sup>, Hirokazu Kameoka<sup>2</sup>, Tomoki Toda<sup>3</sup>, and Satoshi Nakamura<sup>1</sup>

<sup>1</sup>Graduate School of Information Science, Nara Institute of Science and Technology, Japan

<sup>2</sup>NTT Communication Science Laboratories, NTT Corporation, Japan

<sup>3</sup>Information Technology Center, Nagoya University, Japan

{ko-t, s-nakamura}@is.naist.jp, kameoka.hirokazu@lab.ntt.co.jp, tomoki@icts.nagoya-u.ac.jp

## ABSTRACT

We have previously proposed a statistical fundamental frequency ( $F_0$ ) prediction method that makes it possible to predict the underlying  $F_0$  contour of electrolaryngeal (EL) speech from its spectral feature sequence. Although this method was shown to contribute to improving the naturalness of EL speech as a whole, the predicted  $F_0$  contour was still unnatural compared with that in normal speech. One possible solution to improve the naturalness of the predicted  $F_0$  contours would be to take account of the physical mechanism of vocal phonation. Recently a statistical model of voice  $F_0$  contours was formulated by constructing a stochastic counterpart of the Fujisaki model, a well-founded mathematical model representing the control mechanism of vocal fold vibration. This paper proposes a Product-of-Experts model to incorporate this generative model of voice  $F_0$  contours into the statistical  $F_0$  prediction model. Based on the constructed model, we derive algorithms for parameter training and  $F_0$  prediction. Experimental results revealed that the proposed method successfully outperformed our previously proposed method in terms of the naturalness of the predicted  $F_0$  contours.

**Index Terms**— Electrolaryngeal speech enhancement,  $F_0$  prediction, Generative model, Product of Experts

## 1. INTRODUCTION

Speech is one of the most common tools in human communication. Since speech is produced by the vocal apparatus, the produced sounds are physically constrained by the conditions of human body. Unfortunately, there are many people with disabilities that prevent them from producing speech freely, leading to communication barriers. Those who are unable to produce speech freely involve laryngectomees, who have undergone an operation to remove the larynx including the vocal folds for such reasons as injury or laryngeal cancer. The ability by these people to generate excitation sounds is severely impaired because they no longer have their vocal folds. One alternative means of producing voice for these patients involves the use of electrolaryngeal (EL) speech, which is produced by using the excitation signals mechanically generated from an electrolarynx. EL speech is reasonably intelligible, but somewhat unnatural particularly due to the mechanical sounding of the excitation signals.

To address this problem, we have previously proposed methods that aim to convert EL speech to normal-sounding speech, by predicting the fundamental frequency ( $F_0$ ) contour from the spectrum sequence of the EL speech based on

Gaussian Mixture Models (GMMs) followed by synthesizing the speech waveforms according to the predicted acoustic parameters [1–3]. These methods were shown to contribute to improving the naturalness of EL speech [1, 2] and also preserving its intelligibility [3]. However, the  $F_0$  contours predicted using these methods still sounded unnatural compared with that in normal speech. This was because the predicted  $F_0$  contours were not necessarily guaranteed to satisfy the physical constraint of the actual control mechanism of the thyroid cartilage, even though they were optimal in a statistical sense. In this regard, these methods still had a plenty of room for improvement. One possible solution to improve the naturalness of the  $F_0$  contours of the converted speech would be to incorporate a generative model of voice  $F_0$  contours into the statistical  $F_0$  prediction model to take account of the physical mechanism of vocal phonation.

One of the authors previously proposed a statistical model of voice  $F_0$  contours [4–6], formulated by constructing a stochastic counterpart of the Fujisaki model [7], a well-founded mathematical model representing the control mechanism of vocal fold vibration. The Fujisaki model [7] assumes that an  $F_0$  contour on a logarithmic scale is the superposition of a phrase component, an accent component and a base value. 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 model proposed in [4–6] has made it possible to estimate the underlying parameters of the Fujisaki model that best explain the given  $F_0$  contour, by using powerful statistical inference techniques.

To incorporate the generative  $F_0$  contour model into the statistical  $F_0$  prediction framework, this paper proposes a Product-of-Experts (PoE) model [8] combining the above-mentioned two models. Since the PoE model is obtained by multiplying the densities of different models, it usually becomes complicated due to the renormalization term. To avoid this, we introduce a latent trajectory model proposed in [9] to reformulate the prediction model so that it can be smoothly combined with the generative  $F_0$  contour model.

## 2. GMM-BASED STATISTICAL $F_0$ PREDICTION

We briefly review our statistical  $F_0$  prediction method [1–3], which exploits the idea of statistical voice conversion techniques [10, 11]. The aim of this method is to predict  $F_0$  contours from the spectral parameters of EL speech. As with voice conversion methods, it consists of training and prediction processes.

In the training process, the parameters  $\lambda_G$  of the joint probability density  $p((\mathbf{x}[k]^T, \mathbf{o}[k]^T)^T | \lambda_G)$  described as a Gaussian mixture model (GMM) are trained, where  $\mathbf{x}$  de-

This work was supported in part of JSPS KAKENHI 26280060.

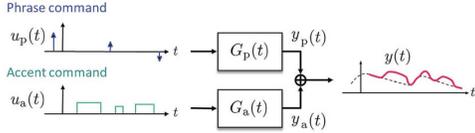


Fig. 1. Original Fujisaki model [7].

notes transposition, and  $\mathbf{x}[k]$  and  $\mathbf{o}[k]$  denote a source feature and a target feature at time frame  $k$ , respectively. The corresponding joint feature vectors can be obtained by performing automatic frame alignment with Dynamic Time Warping. As a source feature, the spectral segment feature of EL speech is extracted on a frame-by-frame basis from the mel-cepstra at multiple frames around the current frame  $k$  [12]. The target feature  $\mathbf{o}[k] = (y[k], \Delta y[k])^T$  consists of the static and delta (time derivative) components of the log-scaled  $F_0$  value  $y[k]$ , extracted on a frame-by-frame basis from the target normal speech. Note that to improve prediction accuracy, we interpolate unvoiced frames of  $F_0$  patterns by using spline interpolation and remove micro-prosody [13].

In the prediction process, given the spectral segment sequence  $\mathbf{x} = (\mathbf{x}[1]^T, \dots, \mathbf{x}[K]^T)^T$  of EL speech, the most likely  $F_0$  sequence  $\mathbf{y} = (y[1], \dots, y[K])^T$  can be obtained as follows:

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} p(\mathbf{o}|\mathbf{x}, \lambda_G) \quad \text{subject to } \mathbf{o} = \mathbf{W}\mathbf{y}, \quad (1)$$

where  $\mathbf{o} = (\mathbf{o}[1]^T, \dots, \mathbf{o}[K]^T)^T$  denotes the joint static and dynamic feature vector sequence,  $\mathbf{W}$  is a constant matrix that transforms the static feature vector sequence  $\mathbf{y}$  to  $\mathbf{o}$ . Namely, each row of  $\mathbf{W}$  consists of the coefficients of an identity mapping operator or time differential operator.  $p(\mathbf{o}|\mathbf{x}, \lambda_G)$  is the GMM with the trained parameters, which we approximate as

$$\begin{aligned} p(\mathbf{o}|\mathbf{x}, \lambda_G) &= \sum_{\mathbf{m}} p(\mathbf{o}|\mathbf{x}, \mathbf{m}, \lambda_G) p(\mathbf{m}|\mathbf{x}, \lambda_G) \\ &\simeq p(\mathbf{o}|\mathbf{x}, \hat{\mathbf{m}}, \lambda_G) p(\hat{\mathbf{m}}|\mathbf{x}, \lambda_G), \end{aligned} \quad (2)$$

with  $\hat{\mathbf{m}} = \operatorname{argmax}_{\mathbf{m}} p(\mathbf{m}|\mathbf{x}, \lambda_G)$ . Here,  $\mathbf{m} = (m_1, \dots, m_K)$  indicates a sequence of mixture indices.  $p(\mathbf{o}|\mathbf{x}, \mathbf{m}, \lambda_G)$  is given as the product of  $p(\mathbf{o}[k]|\mathbf{x}[k], m_k, \lambda_G) = \mathcal{N}(\mathbf{o}[k]; \mathbf{e}_{m_k}^{(\mathbf{o}|\mathbf{x})}, \mathbf{D}_{m_k}^{(\mathbf{o}|\mathbf{x})})$  over  $k$  where  $\mathbf{e}_{m_k}^{(\mathbf{o}|\mathbf{x})}$  and  $\mathbf{D}_{m_k}^{(\mathbf{o}|\mathbf{x})}$  are the mean vector and covariance matrix of  $m_k$ -th mixture component, respectively. Thanks to this approximation, the solution to Eq. (1) is given explicitly as follows:

$$\hat{\mathbf{y}} = (\mathbf{W}^T \mathbf{D}^{(\mathbf{o}|\mathbf{x})^{-1}} \mathbf{W})^{-1} \mathbf{W}^T \mathbf{D}^{(\mathbf{o}|\mathbf{x})^{-1}} \mathbf{e}_{\hat{\mathbf{m}}}^{(\mathbf{o}|\mathbf{x})}, \quad (3)$$

where  $\mathbf{e}_{\hat{\mathbf{m}}}^{(\mathbf{o}|\mathbf{x})}$  is a stacked vector of the mean vectors  $\mathbf{e}_{\hat{m}_1}^{(\mathbf{o}|\mathbf{x})}, \dots, \mathbf{e}_{\hat{m}_K}^{(\mathbf{o}|\mathbf{x})}$  and  $\mathbf{D}_{\hat{\mathbf{m}}}^{(\mathbf{o}|\mathbf{x})}$  is a block diagonal matrix where each block is  $\mathbf{D}_{\hat{m}_1}^{(\mathbf{o}|\mathbf{x})}, \dots, \mathbf{D}_{\hat{m}_K}^{(\mathbf{o}|\mathbf{x})}$ .

### 3. GENERATIVE MODEL OF VOICE $F_0$ CONTOURS

The generative model of  $F_0$  contours proposed in [4–6] is a stochastic counterpart of a discrete-time version of the Fujisaki model [7].

The Fujisaki model (shown in Fig. 1) assumes that a log-scaled  $F_0$  contour  $y(t)$  is the superposition of a phrase component  $y_p(t)$ , an accent component  $y_a(t)$  and a base value  $\mu_b$ . The phrase and accent components are assumed to be the outputs of different second-order critically damped filters, excited with Dirac deltas  $u_p(t)$  (phrase commands) and rectangular pulses  $u_a(t)$  (accent commands), respectively. Here,

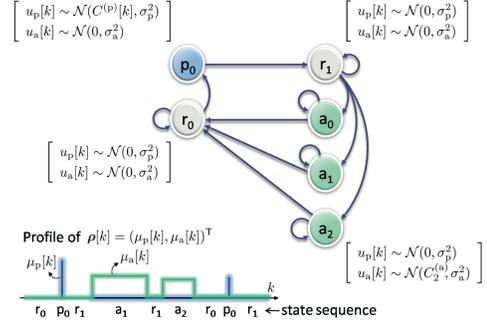


Fig. 2. Command function modeling with HMM.

it must be noted that the phrase and accent commands do not usually overlap each other. The base value 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

$$y(t) = y_p(t) + y_a(t) + \mu_b, \quad (4)$$

where

$$y_p(t) = g_p(t) * u_p(t), \quad (5)$$

$$y_a(t) = g_a(t) * u_a(t). \quad (6)$$

Here,  $*$  denotes convolution over time.  $g_p(t)$  and  $g_a(t)$  are the impulse responses 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.

A key idea in the proposed model [4–6] is that the sequence of the phrase and accent command pair (i.e., the underlying parameters of the Fujisaki model) is modeled as a path-restricted hidden Markov model (HMM) with Gaussian emission densities (shown in Fig. 2) so that estimating the state transition of the HMM directly amounts to estimating the Fujisaki-model parameters.

We hereafter use  $k$  to indicate the discrete time index. Given a state sequence  $\mathbf{s} = (s_1, \dots, s_K)$  of the above HMM, the conditional distributions of the phrase command sequence  $\mathbf{u}_p = (u_p[1], \dots, u_p[K])^T$  and the accent command sequence  $\mathbf{u}_a = (u_a[1], \dots, u_a[K])^T$  are given as

$$p(\mathbf{u}_p|\mathbf{s}, \lambda_F) = \mathcal{N}(\mathbf{u}_p; \boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p), \quad (7)$$

$$p(\mathbf{u}_a|\mathbf{s}, \lambda_F) = \mathcal{N}(\mathbf{u}_a; \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a), \quad (8)$$

respectively, where  $\lambda_F$  denotes the parameters of the HMM.  $\boldsymbol{\mu}_p$  and  $\boldsymbol{\mu}_a$  denote the mean sequences of the state emission densities and  $\boldsymbol{\Sigma}_p$  and  $\boldsymbol{\Sigma}_a$  are diagonal matrices whose diagonal elements correspond to the variances of the state emission densities. From Eqs. (5) and (6), the relationships between  $\mathbf{y}_p = (y_p[1], \dots, y_p[K])^T$  and  $\mathbf{u}_p$  and between  $\mathbf{y}_a = (y_a[1], \dots, y_a[K])^T$  and  $\mathbf{u}_a$  can be written as

$$\mathbf{G}_p \mathbf{u}_p = \mathbf{y}_p, \quad (9)$$

$$\mathbf{G}_a \mathbf{u}_a = \mathbf{y}_a, \quad (10)$$

where  $\mathbf{G}_p$  and  $\mathbf{G}_a$  are Toeplitz matrices where each row is a shifted copy of the convolution kernels  $g_p[1], \dots, g_p[K]$  and  $g_a[1], \dots, g_a[K]$ . By using  $\mathbf{u}_b$  to denote the baseline component, the log  $F_0$  sequence  $\mathbf{y}$  is given as  $\mathbf{y} = \mathbf{y}_p + \mathbf{y}_a + \mathbf{u}_b + \mathbf{n}$  where  $\mathbf{n}$  is an additive noise component corresponding to micro prosody. If we assume that  $\mathbf{n}$  follows a Gaussian distribution with mean  $\mathbf{0}$  and covariance  $\mathbf{\Gamma}$ , the conditional distribution of  $\mathbf{y}$  given  $\mathbf{u} = (\mathbf{u}_p^T, \mathbf{u}_a^T, \mathbf{u}_b^T)^T$  is defined as

$$p(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}; \mathbf{G}_p \mathbf{u}_p + \mathbf{G}_a \mathbf{u}_a + \mathbf{u}_b, \mathbf{\Gamma}). \quad (11)$$

We further assume that  $\mathbf{u}_b$  follows a Gaussian distribution with mean  $\mu_b \mathbf{1}$  and covariance  $\mathbf{\Sigma}_b$ . Then, from Eqs. (7), (8) and (11), the conditional distribution of  $\mathbf{y}$  given  $\mathbf{s}$  is given as

$$\begin{aligned} p(\mathbf{y}|\mathbf{s}, \mathbf{\Lambda}_F) &= \int p(\mathbf{y}|\mathbf{u})p(\mathbf{u}|\mathbf{s}, \mathbf{\lambda}_F) d\mathbf{u} \\ &= \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_F, \mathbf{\Sigma}_F), \end{aligned} \quad (12)$$

where  $\boldsymbol{\mu}_F = \mathbf{G}_p \boldsymbol{\mu}_p + \mathbf{G}_a \boldsymbol{\mu}_a + \mu_b \mathbf{1}$  and  $\mathbf{\Sigma}_F = \mathbf{G}_p \mathbf{\Sigma}_p \mathbf{G}_p^\top + \mathbf{G}_a \mathbf{\Sigma}_a \mathbf{G}_a^\top + \mathbf{\Sigma}_b + \mathbf{\Gamma}$ .

## 4. PROPOSED MODEL

### 4.1. Product-of-Experts Strategy

PoE [8] is a general technique to model a complicated distribution of data by combining relatively simpler distributions (experts). Since the distribution is obtained by multiplying the densities of the experts, the way the experts are combined is somewhat similar to an ‘‘and’’ operation. In this section, we construct a PoE model by treating the two models introduced in the previous sections as the experts.

Training a PoE model by maximizing the likelihood of the data usually becomes difficult since it is hard even to approximate the derivatives of the renormalization term. By contrast, we propose an elegant formulation that allows the use of the EM algorithm for both parameter training and  $F_0$  prediction. To do so, we first introduce a latent trajectory model proposed in [9] to reformulate the GMM-based statistical  $F_0$  prediction model, which plays a key role in making this possible.

### 4.2. Latent-Trajectory-GMM-Based $F_0$ Prediction

We reformulate the GMM-based statistical  $F_0$  prediction model presented in Sec. 2 by employing the idea proposed in [9]. Instead of treating  $\mathbf{o}$  as a function of  $\mathbf{y}$ , we treat  $\mathbf{o}$  as a latent variable to be marginalized out, that is related to  $\mathbf{y}$  through a soft constraint  $\mathbf{o} \simeq \mathbf{W}\mathbf{y}$ . The relationship  $\mathbf{o} \simeq \mathbf{W}\mathbf{y}$  can be expressed through the conditional distribution  $p(\mathbf{y}|\mathbf{o})$

$$\begin{aligned} p(\mathbf{y}|\mathbf{o}) &\propto \exp \left\{ -\frac{1}{2} (\mathbf{W}\mathbf{y} - \mathbf{o})^\top \mathbf{\Lambda} (\mathbf{W}\mathbf{y} - \mathbf{o}) \right\} \\ &= \mathcal{N}(\mathbf{y}; \mathbf{H}\mathbf{o}, \mathbf{V}), \end{aligned} \quad (13)$$

where  $\mathbf{H} = (\mathbf{W}^\top \mathbf{\Lambda} \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{\Lambda}$  and  $\mathbf{V} = (\mathbf{W}^\top \mathbf{\Lambda} \mathbf{W})^{-1}$ .  $\mathbf{\Lambda}$  is a constant positive definite matrix that can be set arbitrarily. As with Sec. 2, the joint distribution  $p(\mathbf{x}, \mathbf{o}|\lambda_G)$  is modeled as a GMM. Namely, given mixture indices  $\mathbf{m}$ , the conditional distribution  $p(\mathbf{x}, \mathbf{o}|\mathbf{m}, \lambda_G)$  is defined as a Gaussian distribution. Thus, the joint distribution of  $\mathbf{y}$ ,  $\mathbf{x}$ ,  $\mathbf{o}$ , and  $\mathbf{m}$  can be described using the distributions defined above

$$p(\mathbf{y}, \mathbf{x}, \mathbf{o}, \mathbf{m}|\lambda_G) = p(\mathbf{y}|\mathbf{o})p(\mathbf{x}, \mathbf{o}|\mathbf{m}, \lambda_G)p(\mathbf{m}|\lambda_G), \quad (15)$$

where  $p(\mathbf{m}|\lambda_G)$  is the product of mixture weights of the GMM. By marginalizing  $\mathbf{o}$  and  $\mathbf{m}$  out, we can readily obtain the joint distribution  $p(\mathbf{y}, \mathbf{x}|\lambda_G)$ , which can be used as a criterion to train  $\lambda_G$  and predict optimal  $\mathbf{y}$  in a consistent manner, unlike the method presented in Sec. 2. We use this model to construct our PoE model in the next subsection.

### 4.3. Deriving PoE

In the same way as Eq. (15), we write the model presented in Sec. 3 in the form of a joint distribution

$$p(\mathbf{y}, \mathbf{u}, \mathbf{s}|\lambda_F) = p(\mathbf{y}|\mathbf{u})p(\mathbf{u}|\mathbf{s}, \lambda_F)p(\mathbf{s}|\lambda_F), \quad (16)$$

where  $p(\mathbf{u}|\mathbf{s}, \lambda_F)$  is given as the product of state emission densities and  $p(\mathbf{s}|\lambda_F)$  the product of the state transition

probabilities given a state sequence  $\mathbf{s}$ . We consider constructing a PoE model by combining Eqs. (15) and (16) followed by marginalization, rather than by simply combining the marginal distributions  $p(\mathbf{y}, \mathbf{x}|\lambda_G)$  and  $p(\mathbf{y}|\lambda_F)$ , which makes the parameter training and  $F_0$  prediction problems excessively hard. To do so, we first combine the densities of  $p(\mathbf{y}|\mathbf{o})$  and  $p(\mathbf{y}|\mathbf{u})$  to obtain  $p(\mathbf{y}|\mathbf{o}, \mathbf{u})$ . Since both of these distributions are Gaussians, the product of their distributions can be easily obtained by completing the square of the exponent

$$\begin{aligned} p(\mathbf{y}|\mathbf{o}, \mathbf{u}) &\propto \mathcal{N}(\mathbf{y}; \mathbf{H}\mathbf{o}, \mathbf{V}) \cdot \mathcal{N}(\mathbf{y}; \mathbf{G}\mathbf{u}, \mathbf{\Gamma}) \\ &= \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_{y|\mathbf{o}, \mathbf{u}}, \mathbf{\Sigma}_{y|\mathbf{o}, \mathbf{u}}), \end{aligned} \quad (17)$$

$$\boldsymbol{\mu}_{y|\mathbf{o}, \mathbf{u}} = (\mathbf{V}^{-1} + \mathbf{\Gamma}^{-1})^{-1} (\mathbf{V}^{-1} \mathbf{H}\mathbf{o} + \mathbf{\Gamma}^{-1} \mathbf{G}\mathbf{u}), \quad (18)$$

$$\mathbf{\Sigma}_{y|\mathbf{o}, \mathbf{u}} = (\mathbf{V}^{-1} + \mathbf{\Gamma}^{-1})^{-1}, \quad (19)$$

where  $\mathbf{G} = [\mathbf{G}_p \mathbf{G}_a \mathbf{I}]$  and  $\mathbf{u} = (\mathbf{u}_p^\top, \mathbf{u}_a^\top, \mathbf{u}_b^\top)^\top$ . From Eqs. (15), (16) and (17), the joint distribution of  $\mathbf{y}$ ,  $\mathbf{x}$ ,  $\mathbf{o}$ ,  $\mathbf{u}$ ,  $\mathbf{m}$  and  $\mathbf{s}$  can be constructed as

$$\begin{aligned} p(\mathbf{y}, \mathbf{x}, \mathbf{o}, \mathbf{u}, \mathbf{m}, \mathbf{s}|\lambda_G, \lambda_F) & \\ = \underbrace{p(\mathbf{y}|\mathbf{o}, \mathbf{u})p(\mathbf{x}, \mathbf{o}|\mathbf{m}, \lambda_G)p(\mathbf{u}|\mathbf{s}, \lambda_F)}_{p(\mathbf{y}, \mathbf{x}, \mathbf{o}, \mathbf{u}|\mathbf{m}, \mathbf{s}, \lambda_G, \lambda_F)} p(\mathbf{m}|\lambda_G)p(\mathbf{s}|\lambda_F). & \end{aligned} \quad (20)$$

This can be used as the complete data likelihood for parameter training and  $F_0$  prediction as explained later. By marginalizing  $\mathbf{o}$  and  $\mathbf{u}$  out, we can readily obtain the joint distribution  $p(\mathbf{y}, \mathbf{x}, \mathbf{m}, \mathbf{s}|\lambda_G, \lambda_F)$ , which can be used as a criterion to train  $\lambda_G$  and  $\lambda_F$  and predict optimal  $\mathbf{y}$  in a consistent manner.

Since both  $p(\mathbf{x}, \mathbf{o}|\mathbf{m}, \lambda_G)$  and  $p(\mathbf{u}|\mathbf{s}, \lambda_F)$  are Gaussians, let us write them as

$$p(\mathbf{x}, \mathbf{o}|\mathbf{m}, \lambda_G) = \mathcal{N} \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{o} \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_o \end{bmatrix}, \begin{bmatrix} \mathbf{P}_{xx} & \mathbf{P}_{xo} \\ \mathbf{P}_{ox} & \mathbf{P}_{oo} \end{bmatrix}^{-1} \right), \quad (21)$$

$$p(\mathbf{u}|\mathbf{s}, \lambda_F) = \mathcal{N}(\mathbf{u}; \boldsymbol{\mu}_u, \mathbf{P}_u^{-1}). \quad (22)$$

Then, from Eqs. (17), (21) and (22), it can be shown that  $p(\mathbf{y}, \mathbf{x}, \mathbf{o}, \mathbf{u}|\mathbf{m}, \mathbf{s}, \lambda_G, \lambda_F)$  is given as

$$\begin{aligned} p(\mathbf{y}, \mathbf{x}, \mathbf{o}, \mathbf{u}|\mathbf{m}, \mathbf{s}, \lambda_G, \lambda_F) & \\ = \mathcal{N} \left( \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \\ \mathbf{o} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} \mathbf{A}_{11} \mathbf{b}_1 + \mathbf{A}_{12} \mathbf{b}_2 \\ \mathbf{A}_{21} \mathbf{b}_1 + \mathbf{A}_{22} \mathbf{b}_2 \end{bmatrix}, \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \right), & \end{aligned} \quad (23)$$

where

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{\Sigma}_{y|\mathbf{o}, \mathbf{u}}^{-1} & \mathbf{O} & -\mathbf{V}^{-1} \mathbf{H} & -\mathbf{\Gamma}^{-1} \mathbf{G} \\ \mathbf{O} & \mathbf{P}_{xx} & -\mathbf{P}_{xo} & \mathbf{O} \\ -\mathbf{H}^\top \mathbf{V}^{-\top} & -\mathbf{P}_{ox} & \mathbf{P}_{oo} & \mathbf{O} \\ -\mathbf{G}^\top \mathbf{\Gamma}^{-\top} & \mathbf{O} & \mathbf{O} & \mathbf{P}_u \end{bmatrix}^{-1}, \quad (24)$$

$$\begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{O} \\ \mathbf{P}_{xx} \boldsymbol{\mu}_x - \mathbf{P}_{xo} \boldsymbol{\mu}_o \\ \mathbf{P}_{oo} \boldsymbol{\mu}_o - \mathbf{P}_{ox} \boldsymbol{\mu}_x \\ \mathbf{P}_u \boldsymbol{\mu}_u \end{bmatrix}, \quad (25)$$

by completing the square of the exponent. Note that  $\mathbf{A}_{11}$ ,  $\mathbf{A}_{12}$ ,  $\mathbf{A}_{21}$  and  $\mathbf{A}_{22}$  can be written explicitly using the block-wise inversion formula.

### 4.4. Parameter Training and $F_0$ Prediction

The problems of parameter training and  $F_0$  prediction can be formulated as the following optimization problems:

$$\{\hat{\lambda}_G, \hat{\lambda}_F, \hat{\mathbf{m}}, \hat{\mathbf{s}}\} = \underset{\lambda_G, \lambda_F, \mathbf{m}, \mathbf{s}}{\operatorname{argmax}} \log p(\hat{\mathbf{y}}, \hat{\mathbf{x}}, \hat{\mathbf{m}}, \hat{\mathbf{s}}|\lambda_G, \lambda_F), \quad (26)$$

$$\{\hat{\mathbf{y}}, \hat{\mathbf{m}}, \hat{\mathbf{s}}\} = \underset{\mathbf{y}, \mathbf{m}, \mathbf{s}}{\operatorname{argmax}} \log p(\mathbf{y}, \hat{\mathbf{x}}, \mathbf{m}, \hat{\mathbf{s}}|\lambda_G, \hat{\lambda}_F). \quad (27)$$

where  $\tilde{\mathbf{y}}$  and  $\tilde{\mathbf{x}}$  denote the observed  $F_0$  contour extracted from normal speech and the observed spectral sequence extracted from non-larynx speech. Both of these problems can be solved using the EM algorithm by treating  $\mathbf{o}$  and  $\mathbf{u}$  as latent variables. Owing to space limitations, here we only derive an algorithm for solving Eq. (27).

The likelihood of  $\mathbf{y}$ ,  $\mathbf{m}$  and  $\mathbf{s}$  given the complete data  $\{\tilde{\mathbf{x}}, \mathbf{o}, \mathbf{u}\}$  is given by Eq. (20). By taking the conditional expectation of  $\log p(\mathbf{y}, \tilde{\mathbf{x}}, \mathbf{o}, \mathbf{u} | \mathbf{m}, \mathbf{s}, \hat{\lambda}_G, \hat{\lambda}_F)$  with respect to  $\mathbf{o}$  and  $\mathbf{u}$  given  $\tilde{\mathbf{x}}, \mathbf{y} = \mathbf{y}', \mathbf{m} = \mathbf{m}'$  and  $\mathbf{s} = \mathbf{s}'$  and then adding  $\log p(\mathbf{m} | \hat{\lambda}_G) p(\mathbf{s} | \hat{\lambda}_F)$ , we obtain an auxiliary function

$$Q(\theta, \theta') = \mathbb{E}_{\mathbf{o}, \mathbf{u} | \tilde{\mathbf{x}}, \mathbf{y}', \mathbf{m}', \mathbf{s}'} [\log p(\mathbf{y}, \tilde{\mathbf{x}}, \mathbf{o}, \mathbf{u} | \mathbf{m}, \mathbf{s}, \hat{\lambda}_G, \hat{\lambda}_F)] + \log p(\mathbf{m} | \hat{\lambda}_G) + \log p(\mathbf{s} | \hat{\lambda}_F), \quad (28)$$

where  $\theta = \{\mathbf{y}, \mathbf{m}, \mathbf{s}\}$ . From Eq. (23), we obtain

$$\mathbb{E} \begin{bmatrix} \mathbf{o} \\ \mathbf{u} \end{bmatrix} \begin{bmatrix} \mathbf{y}' \\ \tilde{\mathbf{x}} \end{bmatrix}, \mathbf{m}', \mathbf{q}' = \mathbf{A}_{21} \mathbf{b}_1 + \mathbf{A}_{22} \mathbf{b}_2 + \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \left( \begin{bmatrix} \mathbf{y}' \\ \tilde{\mathbf{x}} \end{bmatrix} - \mathbf{A}_{11} \mathbf{b}_1 - \mathbf{A}_{12} \mathbf{b}_2 \right) =: \begin{bmatrix} \bar{\mathbf{o}} \\ \bar{\mathbf{u}} \end{bmatrix}, \quad (29)$$

$$\mathbb{E} \begin{bmatrix} \mathbf{o} \\ \mathbf{u} \end{bmatrix} \begin{bmatrix} \mathbf{o} \\ \mathbf{u} \end{bmatrix}^T \begin{bmatrix} \mathbf{y}' \\ \tilde{\mathbf{x}} \end{bmatrix}, \mathbf{m}', \mathbf{q}' = \mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12} + \begin{bmatrix} \bar{\mathbf{o}} \\ \bar{\mathbf{u}} \end{bmatrix} \begin{bmatrix} \bar{\mathbf{o}} \\ \bar{\mathbf{u}} \end{bmatrix}^T, \quad (30)$$

which are the values to be computed at the ‘‘E-step’’ by substituting  $\theta$  into  $\theta'$ . At the ‘‘M-step’’, we compute

$$\{\mathbf{y}, \mathbf{m}, \mathbf{s}\} \leftarrow \underset{\mathbf{y}, \mathbf{m}, \mathbf{s}}{\operatorname{argmax}} Q(\theta, \theta'). \quad (31)$$

The update equations are omitted owing to space limitations.

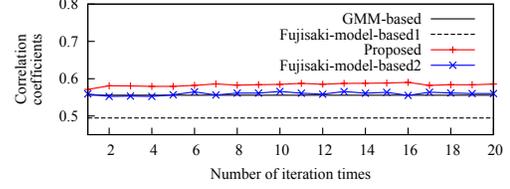
## 5. EXPERIMENTAL EVALUATION

### 5.1. Experimental Conditions

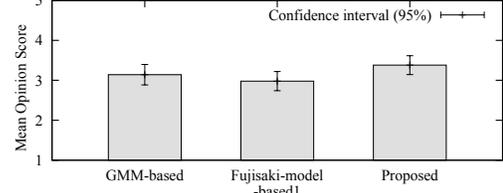
We conducted objective and subjective evaluation experiments to evaluate the performance of the proposed method. For the objective evaluation, we evaluated the  $F_0$  correlation coefficients between the predicted and target  $F_0$  contours. We also subjectively evaluated the naturalness of the  $F_0$  contour of converted speech.

The source speech was EL speech uttered by one male laryngectomee, and the target speech was normal speech uttered by a professional female speaker. Each speaker uttered about 50 sentences in the ATR phonetically balanced sentence set [15]. We conducted a 5-fold cross validation test in which 40 utterance pairs were used for training, and the remaining 10 utterance pairs were used for evaluation. The sampling frequency was set at 16 kHz. The settings of HMM in the generative  $F_0$  contour model were the same as reported in [4–6].

For initialization of the mixture index sequence  $\mathbf{m}$  and state sequence  $\mathbf{s}$ , we performed the conventional GMM-based and Fujisaki-model-based methods. Note that the Fujisaki-model-based method refers to a post-processing method that consists of first applying the GMM-based method and then fitting the Fujisaki model to the predicted  $F_0$  contour using the method of [5, 6], which is similar to a post-processing method for HMM-based speech synthesis [16]. Note that in this experiment, we implemented a simplified and approximated version of the proposed method, in which the E-step procedure is replaced with the conventional Fujisaki-model-based and GMM-based methods. Therefore, the convergence of the algorithm implemented for the current evaluation is not strictly guaranteed. The speech used for evaluation were synthesized using STRAIGHT [17] given the mel-cepstrum sequence and  $F_0$  contour. The methods selected for comparison were:



**Fig. 3.**  $F_0$  correlation coefficients between predicted  $F_0$  patterns from EL speech and extracted  $F_0$  patterns from normal speech.



**Fig. 4.** Result of opinion test on naturalness.

- *GMM-based*: Predict  $F_0$  contours with the GMM-based method.
- *Fujisaki-model-based1*: Fit the Fujisaki model to the predicted  $F_0$  contours obtained with the *GMM-based* method.
- *Proposed*: Predict  $F_0$  contours with an approximated version of the proposed method, in which the E-step is replaced with the *GMM-based* and *Fujisaki-model-based2* methods.
- *Fujisaki-model-based2*: Fit the Fujisaki model to the predicted  $F_0$  contours obtained with *Proposed*.

### 5.2. Experimental Results

As Fig. 3 shows, *Proposed* obtained the highest prediction accuracy because of eq. (18) meaning an ‘‘and’’ operation for eq. (11) and (14). Therefore, we found that it is effectiveness to construct a PoE model. Furthermore, since *Fujisaki-model-based2* has higher correlation coefficients than *Fujisaki-model-based1*, we found that the predicted  $F_0$  contours by *Proposed* were given good influences by considering not only the GMM-based method but also the Fujisaki model. Note that we used the predicted  $F_0$  contours obtained with the *GMM-based* method as the input for the *Fujisaki-model-based1* method in this experiment. To make a more fair comparison, it would be necessary to modify the Fujisaki-model-based method so as not to depend on the GMM-based method. In addition, as for *Proposed*, there is no large difference of correlation coefficients in each iteration. To make exactly evaluation, we have to conduct the PoE model not replaced E-step with the conventional methods.

As Fig. 4 shows, *Proposed* outperformed the conventional methods, *GMM-based* and *Fujisaki-model-based1*. This result is reasonable since *Proposed* obtained the highest prediction accuracy as in Fig. 3.

## 6. CONCLUSIONS

In this paper, to improve  $F_0$  prediction performance in electrolaryngeal speech enhancement, we proposed a Product-of-Experts model that combined two conventional methods, a statistical  $F_0$  prediction method and a statistical  $F_0$  contour modeling method based on its generative process. Experimental results revealed that the proposed method successfully outperformed our previously proposed method in terms of the naturalness of the predicted  $F_0$  contours.

## 7. REFERENCES

- [1] K. Nakamura, T. Toda, H. Saruwatari, and K. Shikano, "Speaking-aid systems using gmm-based voice conversion for electrolaryngeal speech," in *Proc. Speech Communication*, vol. 54, pp. 134–146, January 2012.
- [2] H. Doi, T. Toda, K. Nakamura, H. Saruwatari, and K. Shikano, "Alaryngeal speech enhancement based on one-to-many eigenvoice conversion," *Audio, Speech, and Language Processing, IEEE/ACM Transactionson*, vol. 22, no. 1, pp. 172–183, January 2014.
- [3] K. Tanaka, T. Toda, G. Neubig, S. Sakti, and S. Nakamura, "A hybrid approach to electrolaryngeal speech enhancement based on noise reduction and statistical excitation generation," *IEICE Transactions on Information and Systems*, vol. E97-D, no. 6, pp. 1429–1437, June 2014.
- [4] H. Kameoka, J. Le Roux, Y. Ohishi, "A statistical model of speech  $F_0$  contours," *ISCA Tutorial and Research Workshop on Statistical And Perceptual Audition*, pp. 43–48, September 2010.
- [5] K. Yoshizato, H. Kameoka, D. Saito, and S. Sagayama, "Statistical approach to fujisaki-model parameter estimation from speech signals and its quantitative evaluation," *Proc. Speech Prosody*, pp. 175–178, May 2012.
- [6] H. Kameoka, K. Yoshizato, T. Ishihara, K. Kadowaki, Y. Ohishi, and K. Kashino, "Generative modeling of voice fundamental frequency contours," *Audio, Speech, and Language Processing, IEEE/ACM Transactionson*, vol. 23, no. 6, pp. 1042–1053, June 2015.
- [7] H. Fujisaki, "A note on the physiological and physical basis for the phrase and accent components in the voice fundamental frequency contour," *Vocal Fold Physiology: Voice Production, Mechanisms and Functions*, pp. 347–355, 1998.
- [8] G. E Hinton, "Training products of experts by minimizing contrastive divergence," *Neural computation*, vol. 14, no. 8, pp. 1771–1800, August 2002.
- [9] H. Kameoka, "Modeling speech parameter sequences with latent trajectory hidden markov model," in *Proc. The 25th IEEE International Workshop on Machine Learning for Signal Processing*, September 2015.
- [10] Y. Stylianou, O. Cappe, and E. Moulines, "Continuous probabilistic transform for voice conversion," *Speech and Audio Processing, IEEE Transactions on*, vol. 6, no. 2, pp. 131–142, March 1998.
- [11] T. Toda, A.W. Black, and K. Tokuda, "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 8, pp. 2222–2235, November 2007.
- [12] T. Toda, M. Nakagiri, and K. Shikano, "Statistical voice conversion techniques for body-conducted unvoiced speech enhancement," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 9, pp. 2505–2517, November 2012.
- [13] K. J. Kohler, "Macro and micro  $F_0$  in the synthesis of intonation," *Papers in Laboratory Phonology I*, pp. 115–138, 1990.
- [14] 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 International Conference on Acoustics, Speech, and Signal Processing*, vol. 1, pp. 509–5012, 2002.
- [15] M. Abe, Y. Sagisaka, T. Umeda, and H. Kuwabara, "Speech database," *ATR Technical Report*, TR-I-0166, September 1990.
- [16] T. Matsuda, K. Hirose, and N. Minematsu, "Applying generation process model constraint to fundamental frequency contours generated by hidden-Markov-model-based speech synthesis," *Acoustical Science and Technology*, Vol. 33, No. 4, pp. 221–228, 2012.
- [17] H. Kawahara, I. Masuda-Katsuse, and A. Cheveigne, "Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based  $F_0$  extraction: Possible role of a repetitive structure in sounds," in *Proc. Speech Communication*, Vol. 27, No. 3-4, pp. 187–207, April 1999.