

## Modeling of Fundamental Frequency Contour of Thai Expressive Speech using Fujisaki's Model and Structural Model

Suphattharachai Chomphan  
Department of Electrical Engineering,  
Faculty of Engineering at Si Racha,  
Kasetsart University, 199 M.6, Tungsukhla,  
Si Racha, Chonburi, 20230, Thailand

---

**Abstract: Problem statement:** In spontaneous speech communication, prosody is an important factor that must be taken into account, since the prosody effects on not only the naturalness but also the intelligibility of speech. Focusing on synthesis of Thai expressive speech, a number of systems has been developed for years. However, the expressive speech with various speaking styles has not been accomplished. To achieve the generation of expressive speech, we need to model the fundamental frequency ( $F_0$ ) contours accurately to preserve the speech prosody to preserve the quality of speech prosody. **Approach:** This study presents a comparison of two successful  $F_0$  models. One approach is based on the Fujisaki's model which has been applied for many tonal and toneless languages. Another one is based on the structural model which has been conducted primarily for Mandarin Chinese. It is based on the assumption that the behavioral characteristics of vocal-fold elongation in vibration could be approximated by those of a simple forced vibrating system. Therefore this approach has been applied to model Thai expressive speech with best-fit function. Our speech database consists of male and female speech and each one contains 4 different speech styles including angry style, sad style, enjoyable style and reading style. Five sentences are used for each speech style and each sentence includes 100 samples. The speech sample in each group is analyzed for an  $F_0$  contour, subsequently a number of Fujisaki's and structural modeling parameters are extracted for each contour. Thereafter, the parameters are used to synthesis the  $F_0$  contour and then the synthesized contour is compared with that of natural speech by calculating RMS error. **Results:** From the experimental analysis, it has been observed that RMS error of each speech style is different from the others for both models. It also reveals that the RMS error of the Fujisaki's model is higher than that of the structural model for all speech styles. In other words, the structural model gives the better fit for modeling of the  $F_0$  contour of the expressive speech than that of the Fujisaki's model. **Conclusion:** From the finding, it is a definite evidence that the structural model is more appropriate than that of the Fujisaki's model for modeling four different speech styles including angry style, sad style and enjoyable style and reading style.

**Key words:** Noised speech, speech analysis, fundamental frequency contour, speech enhancement, speech synthesis, speaker-independent, local-accent, model parameters, Fujisaki's model

---

### INTRODUCTION

The fundamental frequency of voice speech is the most important feature among all of the features known to carry prosodic information which is an inherently supra-segmental feature of human speech. The  $F_0$  contours of an utterance convey the stress, intonation and rhythmic structures, which determine the naturalness and intelligibility of synthetic speech. As a result, the appropriate modeling of  $F_0$  contour plays a significant role in the speech processing area, e.g., speech recognition, speech synthesis, speech analysis

and speech coding. A number of modeling techniques in the former studies have been performed in various levels of speech units, e.g., utterance level (Saito and Sakamoto, 2002; Li *et al.*, 2004; Tao *et al.*, 2006), word and syllable levels (Fujisaki and Sudo, 1971; Tran *et al.*, 2006). In Thai speech, Fujisaki's model has been successfully applied for modeling of utterances, tones and words (Hiroya and Sumio, 2002; Seresangtakul and Takara, 2002; 2003). In the Thai speech synthesis, the statistical modeling of  $F_0$  contour has been conducted by Chomphan and Kobayashi in the implementation of both speaker-dependent and speaker-

independent systems in 2007-2009 (Chomphan and Kobayashi, 2007; 2008; 2009). Lately, the Fujisaki's model has been applied within a speaker-independent system as extended modules. Moreover, it has also been exploited in the modeling of Thai expressive speech; i.e., sad, happy, angry styles (Chomphan and Kobayashi, 2008; 2009). Another study has been conducted by using a structural model which is based on the assumption that the behavioral characteristics of vocal-fold elongation in vibration could be approximated by those of a simple forced vibrating system (Ni and Hirose, 2006; Chomphan and Kobayashi, 2009). The RMS error calculation has been done for evaluation the modeling performance for both mentioned speech models and also for all speech styles including angry style, sad style, enjoyable style and reading style. This study mainly aims at comparing the Fujisaki's model and the structural model.

### MATERIALS AND METHODS

**Fujisaki's model:** The F0 contour is treated as a linear superposition of a global phrase and local-accent components on a logarithmic scale, as depicted in Fig. 1 (Chomphan and Kobayashi, 2008; 2009).

Typically, the phrase command produces a baseline component, while the accent command produces the accent component of an F0 contour. Mathematically, the F0 contour of an utterance generated from an extension of the Fujisaki's model for tonal languages has the following expressions (Seresangtakul and Takara, 2003):

$$\ln F_0(t) = \ln F_b + \sum_{i=1}^I A_{pi} [G_{pi}(t - T_{0i})] + \sum_{j=1}^J \sum_{k=1}^{K(j)} A_{t,jk} [G_{t,jk}(t - T_{1jk}) - G_{t,jk}(t - T_{2jk})] \quad (1)$$

$$G_{pi}(t) = \begin{cases} (\alpha_i^2 t) \exp(-\alpha_i t) & \text{for } t \geq 0 \\ 0 & \text{for } t < 0 \end{cases} \quad (2)$$

$$G_{t,jk}(t) = \begin{cases} [1 - (1 + \beta_{jk} t) \exp(-\beta_{jk} t)] & \text{for } t \geq 0 \\ 0 & \text{for } t < 0 \end{cases} \quad (3)$$

Where:

$G_{pi}(t)$  = Represents the impulse-response function of the phrase-control mechanism

$G_{t,jk}(t)$  = Represents the step-response function of the tone-control mechanism

The symbols used in three equations denote that  $F_b$  is the smallest F0 value in the F0 contour of interest and  $A_{pi}$  and  $A_{t,jk}$  are the amplitudes of the i-th phrases and of the j-th tone command. Here,  $T_{0i}$  is the timing of the i-th phrase command and  $T_{1jk}$  and  $T_{2jk}$  are the onset and offset of the k-th component of the j-th tone command.  $\alpha_i$  and  $\beta_{jk}$  are time constant parameters, while I, J, K(j) correspond to the number of phrases, tones and components of the j-th tone contained in the utterance.

To find the optimal representative parameters, optimization is carried out by minimizing the mean squared error in the Ln F0(t) domain through the hill-climbing search in the space of model parameters (Seresangtakul and Takara, 2003). To use this model, the parameters are extracted from the speech database, utterance by utterance. The derived parameters are subsequently computed.

**Derived parameters:** From the conventional parameters, we proposed seven derived parameters which reflect the geometrical appearance of the F0 contour of an utterance as follows:

- Baseline frequency
- Numbers of phrase commands
- Numbers of tone commands
- Phrase command duration
- Tone command duration
- Amplitude of phrase command
- Amplitude of tone command

All of them have been extracted for four speech expressions of angry style, sad style, enjoyable style and reading style. Thereafter, the extracted parameters are used to resynthesize the F0 contour in the evaluation process.

**Structural model:** The voice F0 contour is modeled in a logarithmic scale, as depicted in Fig. 2. The mathematical model has been applied (Ni and Hirose, 2006; Chomphan and Kobayashi, 2009) by using a structural control consisting of placing a series of normalized F0 targets along the time axis, which are also specified by transition time and amplitudes. The transitions between targets are approximated by connecting truncated second-order transition functions.

From the background knowledge that the physical factors to regulate the frequency of vocal-fold vibrations are the mass, length and tension of vibrating structures, all of which are dynamically controlled primarily by the intrinsic and extrinsic muscles of the larynx and secondly by the sub-glottal pressure (Ni and Hirose, 2006).

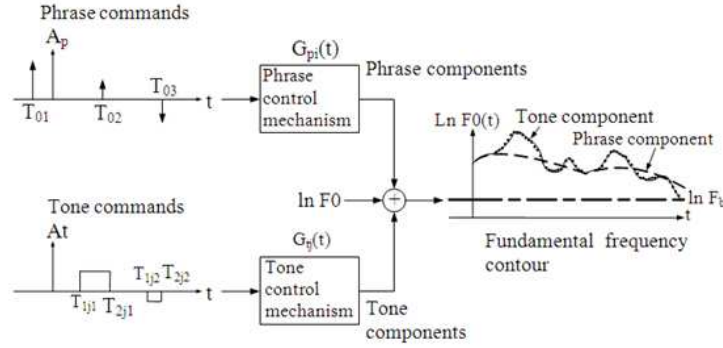


Fig. 1: An extension of Fujisaki's model for the generation of F0 contour

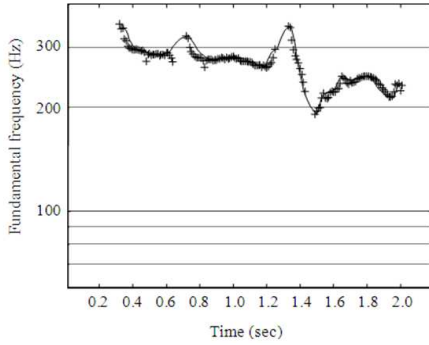


Fig. 2: F<sub>0</sub> contour with a trend line in a logarithmic scale

Fujisaki explained that logarithmic fundamental frequency varies linearly with vocal-fold elongation  $x$  (Fujisaki, 1983), which can be represented in the following mathematical term:

$$\ln f_0 = \frac{b}{2}x + \ln(\sqrt{ac_0}) \quad (4)$$

where,  $a$ ,  $b$  and  $c_0$  are constant coefficients (Fujisaki, 1983).

**Assumption:** The behavioral characteristics of vocal-fold elongation in vibration can be approximated by those of a simple forced vibrating system (Ni and Hirose, 2006).

Formulating the assumption, the behavioral characteristics of a simple forced vibrating system can be characterized by the amplifying coefficients of its vibrating amplitudes:

$$A\left(\frac{\omega_d^2}{\omega_f^2}, \zeta\right) = \frac{1}{\sqrt{\left(1 - \left(1 - 2\zeta^2\right)\frac{\omega_d^2}{\omega_f^2}\right)^2 + 4\zeta^2\left(1 - 2\zeta^2\right)\frac{\omega_d^2}{\omega_f^2}}} \quad (5)$$

where,  $\omega_d$  and  $\omega_f$  denote the natural angular frequencies of the driven system and the driving force, respectively,  $\zeta$  is called damping ratio indicating how tightly the driving force and the driven system are coupled together. Subsequently, replacing  $\omega_d^2 / \omega_f^2$  (square frequency ratio) by  $\lambda$  and substituting  $A(\lambda, \zeta)$  expressed in Eq. 5 for  $x$  of Eq. 4, as a result, the logarithmic fundamental frequency can be expressed as:

$$\ln f_0 = \frac{b \times C}{2} A(\lambda, \zeta) + \ln(\sqrt{ac_0}) \quad (6)$$

where,  $C$  is a constant coefficient.

Typically, a speaker has an individual vocal range. Let  $f_{0t}$  and  $f_{0b}$  denote the top and bottom frequencies of the vocal range of a speaker and  $\lambda_t$  and  $\lambda_b$  denote two  $\lambda$  values that are one-to-one mapped to  $f_{0t}$  and  $f_{0b}$ . The relationship between  $f_0$  within the vocal-range frequency interval and its corresponding  $\lambda$  is shown as follows:

$$\frac{\ln f_0 - \ln f_{0b}}{\ln f_{0t} - \ln f_{0b}} = \frac{A(\lambda, \zeta) - A(\lambda_b, \zeta)}{A(\lambda_t, \zeta) - A(\lambda_b, \zeta)} \quad (7)$$

Since  $f_{0t}$  and  $f_{0b}$  are the top and bottom frequencies of the vocal range,  $\lambda_t$  and  $\lambda_b$  shall be determined regardless of  $\zeta$ .

Practically,  $f_0$  and  $\lambda$  can be determined through  $f_0 = T_{f_0}(\lambda, \zeta)$  and  $\lambda = T_\lambda(f_0, \zeta)$ , where they can be derived from Eq. 7 as followings:

$$T_{f_0}(\lambda, \zeta) = \exp\left(\frac{A(\lambda, \zeta) - A(\lambda_b, \zeta)}{A(\lambda_t, \zeta) - A(\lambda_b, \zeta)} \times \ln \frac{f_{0t}}{f_{0b}} + \ln f_{0b}\right) \quad (8)$$

Table 1: Target values for sparsely specifying an F<sub>0</sub> contour

i	0	1	2	3	4	5	6	7
t <sub>i</sub> (s)	0.32	0.46	0.57	0.72	0.89	1.00	1.18	1.25
f <sub>0i</sub> (Hz)	351.00	296.00	287.00	323.00	275.00	280.00	266.00	296.00
λ <sub>i</sub>	1.28	1.40	1.42	1.34	1.44	1.42	1.46	1.40

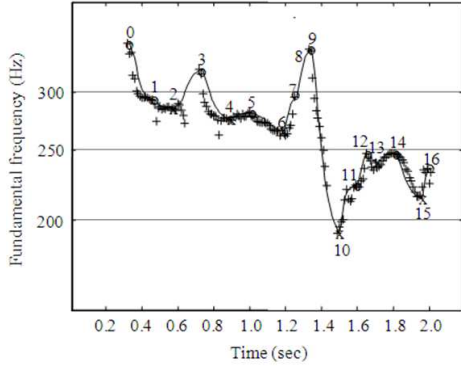


Fig. 3: Example of pitch target allocation on an F<sub>0</sub> contour

and  $T_\lambda(f_0, \zeta)$  can be obtained by searching  $\lambda$  from 1 step-by-step in small increments (e.g., 0.0001), given  $\lambda_b > \lambda_t$ , until  $\lambda$  satisfies the following conditions:

$$\ln \frac{f_0}{f_{0b}} > \frac{A(\lambda + 0.0001, \zeta) - A(\lambda_b, \zeta)}{A(\lambda_t, \zeta) - A(\lambda_b, \zeta)} \times \ln \frac{f_0}{f_{0b}} \quad (9)$$

and

$$\ln \frac{f_0}{f_{0b}} < \frac{A(\lambda - 0.0001, \zeta) - A(\lambda_b, \zeta)}{A(\lambda_t, \zeta) - A(\lambda_b, \zeta)} \times \ln \frac{f_0}{f_{0b}} \quad (10)$$

**Allocation of pitch targets:** By applying the structural model, the pitch targets on an F<sub>0</sub> contour is allocated in advance. Thereafter the parameter of the structural model will be approximated. For example, Fig. 3 shows a sparser specification of the F<sub>0</sub> contours shown in Fig. 1, the ‘o’ signs indicate the tonal peaks, the ‘x’ signs indicate the tonal valleys and the square sign indicates a neutral target to reset an into national phrase.

Let t<sub>i</sub> and f<sub>0i</sub> denote the timing and F<sub>0</sub> value of the i<sup>th</sup> target, respectively. Table 1 lists the first eight target points (t<sub>i</sub>, f<sub>0i</sub>) and corresponding λ<sub>i</sub> = T<sub>λ</sub>(f<sub>0i</sub>, ζ<sub>0</sub>), given λ<sub>t</sub> = 1; λ<sub>b</sub> = 2; f<sub>0b</sub> = 120Hz; f<sub>0t</sub> = 420Hz and ζ<sub>0</sub> = 0.156. The i<sup>th</sup> local F<sub>0</sub> movement, i = 0, ..., 15, is defined as a scope extending from the i<sup>th</sup> target (t<sub>i</sub>, f<sub>0i</sub>) to the next. If f<sub>0i</sub> ≤ f<sub>0i+1</sub>, the local F<sub>0</sub> movement is therefore rising; otherwise, falling. The measured F<sub>0</sub> contours are

subsequently approximated by connecting the rising and falling transitions through these target points.

For an F<sub>0</sub> falling movement, say i = 0, first compute λ(t) for t<sub>0</sub> ≤ t ≤ t<sub>1</sub> by using following equation:

$$\lambda(t) = \lambda_p + \Delta\lambda \left( 1 - \left( 1 + \frac{4.8}{\Delta t} t \right) e^{-\frac{4.8}{\Delta t} t} \right), \quad \text{for } t \geq 0 \quad (11)$$

with the following parameters (based on Table 1): λ<sub>p</sub>(= λ<sub>0</sub>) = 1.28; Δλ = Δλ<sub>0</sub>/0.95 and Δt = Δt<sub>0</sub>/0.95, where Δλ<sub>0</sub> = (Δλ<sub>1</sub> - λ<sub>0</sub>) = 0.12; Δt<sub>0</sub> = (t<sub>1</sub> - t<sub>0</sub>) = 0.14. Second, synthesize contour by using f<sub>0</sub>(t) = T<sub>f0</sub>(λ(t), ζ<sub>0</sub>) of Eq. 5. The thick line between the 0<sup>th</sup> and 1<sup>st</sup> targets shown in Fig. 3 indicates the re-synthesized F<sub>0</sub> contour.

For an F<sub>0</sub> rising movement, say i = 2, first compute λ(t) for t<sub>2</sub> ≤ t ≤ t<sub>3</sub> by using Eq. 8 with the following parameters: λ<sub>p</sub> = 2 - λ<sub>2</sub> = 0.58; Δλ = Δλ<sub>2</sub>/0.95, where Δλ = (2 - λ<sub>3</sub>) - (2 - λ<sub>2</sub>) = 0.08 and Δt = Δt<sub>2</sub>/0.95, where Δt<sub>2</sub> = t<sub>3</sub> - t<sub>2</sub> = 0.15. Then, synthesize contour by using f<sub>0</sub>(t) = T<sub>f0</sub>(2 - λ(t), ζ<sub>0</sub>) of Eq. 5. It has been noted that A(λ, ζ) = A(2 - λ, ζ) is applied to the computation. In Fig. 3, the thick line between the 2<sup>nd</sup> and 3<sup>rd</sup> targets indicates the re-synthesized F<sub>0</sub> contour.

**A mathematical model:** Subsequently, let F<sub>0</sub>(t) represent an F<sub>0</sub> contour as a function of time t in a vocal range [f<sub>0b</sub>, f<sub>0t</sub>] (Ni and Hirose, 2006; Nicknam *et al.*, 2009). Assume Λ(t) to indicate a sequence of virtual tone graphs in λ-time space to specify the underlying lexical tone structures. Additionally, assume a latent scale ζ(t) to characterize the intonation components. Thus, the F<sub>0</sub> contour on the logarithmic scale of fundamental frequency is expressed as a scale transformation from Λ(t) to F<sub>0</sub>(t), corresponding to the syllabic tones fitting themselves with sentence intonation in the vocal range:

$$\frac{\ln F_0(t) - \ln f_{0b}}{\ln f_{0t} - \ln f_{0b}} = \frac{A(\Lambda(t), \zeta(t)) - A(\lambda_b, \zeta(t))}{A(\lambda_t, \zeta(t)) - A(\lambda_b, \zeta(t))} \quad (12)$$

, for t ≥ 0

Where:

$$A(\lambda, \zeta) = \frac{1}{\sqrt{(1 - (1 - 2\zeta^2)\lambda)^2 + 4\zeta^2(1 - 2\zeta^2)\lambda}}, \quad (13)$$

for λ ≥ 1

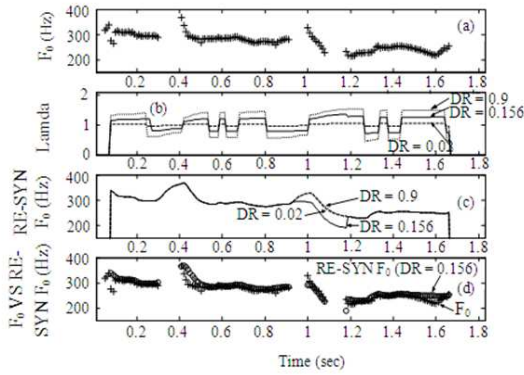


Fig. 4: An example of re-synthesis of  $F_0$  contour by using the structural model (RE-SYN denotes “resynthesized”, DR denotes “damping ratio”)

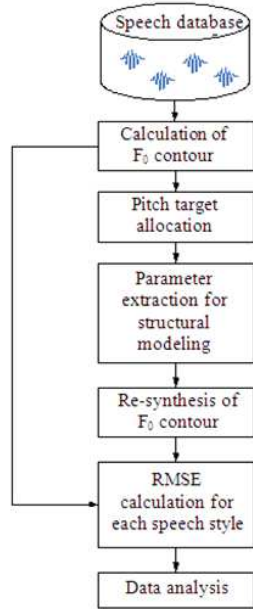


Fig. 5: Work flow for the experimental process

Equation 12-13 jointly indicate a structural formulation of the control process of coupling the syllabic tones and sentence intonation together to form a final sentence melody. Equation 12 states that  $F_0$  contour  $F_0(t)$  is a transformation of a sequence of virtual tone graphs  $\Lambda(t)$  on a latent scale  $\zeta(t)$ .  $\Lambda(t)$  is expressed as a concatenation of  $n$  parametric bell-shaped patterns lining up in series on the time axis with the following definition:

$$\Lambda(t) = \Lambda_{r_i}(t) + \sum_{i=1}^{n-1} \min(\Lambda_{f_i}(t), \Lambda_{r_{i+1}}(t)) + \Lambda_{f_n}(t) \quad (14)$$

where,  $\Lambda_{r_i}(t)$  and  $\Lambda_{f_i}(t)$  indicate the rising and falling transitions of the  $i^{\text{th}}$  bell-shaped pattern, respectively. Their definitions are as follows:

$$\Lambda_{r_i}(t) = \begin{cases} \lambda_{p_i} + \Delta\lambda_{r_i}(1 - D_{r_i}(t_{p_i} - t)), & \text{for } t_{p_{i-1}} \leq t < t_{p_i} \\ 0, & \text{otherwise,} \end{cases} \quad (15)$$

$$\Lambda_{f_i}(t) = \begin{cases} \lambda_{p_i} + \Delta\lambda_{f_i}(1 - D_{f_i}(t - t_{p_i})), & \text{for } t_{p_i} \leq t < t_{p_{i+1}} \\ 0, & \text{otherwise,} \end{cases} \quad (16)$$

Where:

$$D_{x_i}(t) = \left(1 + \frac{4.8t}{\Delta t_{x_i}}\right) e^{-\frac{4.8t}{\Delta t_{x_i}}}, \quad x \in \{r, f\} \quad (17)$$

The model parameters in Eq. 12-17 are defined as follows:

- $[f_0, f_0]$ : Bottom and top frequencies of the vocal Range in hertz
- $[\lambda_b, \lambda_t]$ : Bottom and top values of normalized vocal ranges in  $\lambda$
- $\zeta(t)$ : Latent scales
- $n$ : Number of the bell-shaped patterns
- $(t_{p_i}, \lambda_{p_i})$ :  $i^{\text{th}}$  peak coordinate in  $\lambda$ -time space;  $t_{p_0} = 0$  and  $t_{p_{n+1}} = \infty$
- $\Delta t_{r_i}$ :  $i^{\text{th}}$  rising transition time
- $\Delta\lambda_{r_i}$ :  $i^{\text{th}}$  rising transition amplitude
- $\Delta t_{f_i}$ :  $i^{\text{th}}$  falling transition time and
- $\Delta\lambda_{f_i}$ :  $i^{\text{th}}$  falling transition amplitude,  $i = 1, \dots, n$

Figure 4 shows an example of re-synthesis of  $F_0$  contour by using the structural model. Figure 4a shows the  $F_0$  contour extracted from the natural speech, while Fig. 4b shows corresponding value of  $\lambda$  for three different fixed damping ratios  $\zeta$  (0.156, 0.02 and 0.9). Figure 4c shows the re-synthesized  $F_0$  contour with the three damping ratios in Fig. 4b, while Fig. 4d compares the  $F_0$  contour extracted from the natural speech and the re-synthesized  $F_0$  contour with limited samples (Geravanchizadeh and Rezaii, 2009).

**An experimental design:** The flow chart in Fig. 5 shows the core process for our experiment. At first the speech corpus has been implemented. There is male and female speech in the corpus. Each of them has four speech styles including happy, sad, angry and reading

styles. Each style consists of 5 sentences with 100 samples of utterances. Therefore our speech corpus contains 4,000 utterances. At the beginning, the  $F_0$  values of an utterance have been calculated and then the pitch targets have been allocated by using local minimum/maximum criteria. In between any two adjacent pitch targets used as fixed points, an exponential function has been approximated to minimize the difference between the approximated function and the  $F_0$  contour. The corresponding parameters from all of the functions along the utterance will be used as its representatives. Subsequently, the resynthesis of  $F_0$  contour from the parameters has been conducted. Thereafter, the RMS error between the natural  $F_0$  contour and the resynthesized  $F_0$  contour has been executed. Finally, we analyzed the summarized data from the previous stages.

### RESULTS

The RMS error calculation process has been performed, thereafter the experimental data can be summarized in the following five bar charts (Fig. 6-10). The averaged RMS errors from five different sentences have been calculated. Each bar chart represents one sentence and contains those of four speech styles including happy, sad, angry and reading styles. Moreover each bar chart contains four lines; two are of the Fujisaki's model including male and female speech, while the other two are of the structural model including male and female speech.

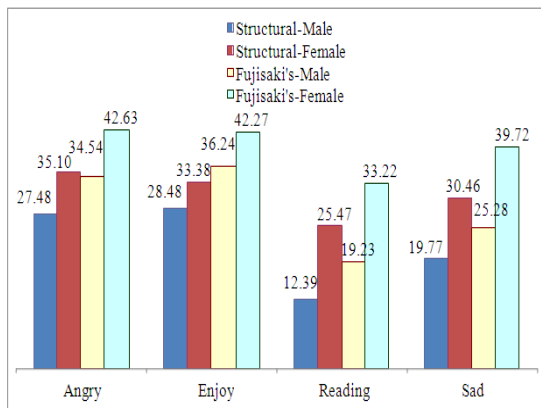


Fig. 6: Averaged RMS error for sentence “κ<sup>η</sup>υν-ταμ-NAv-σεεδ-ρ) & | -φαN” in IPA “have you finished your work?” in English

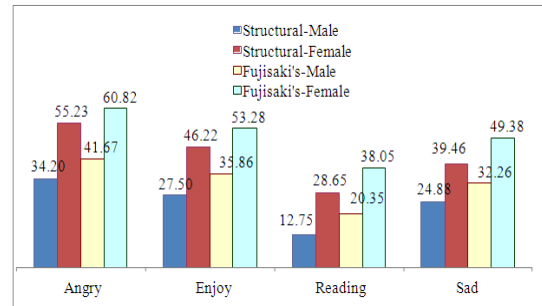


Fig. 7: Averaged RMS error for sentence “τΗΦ | -παφ-να&φ-μα |” in IPA “where have you been?” in English

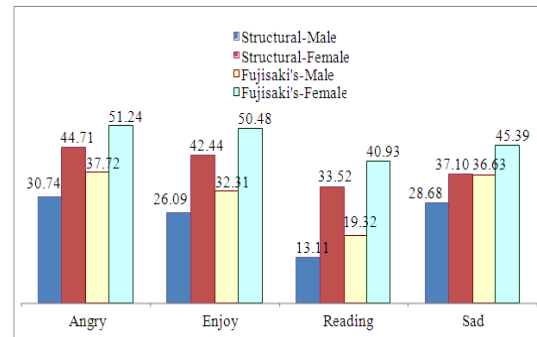


Fig. 8: Averaged RMS error for sentence “χ<sup>η</sup>α&v-χ<sup>η</sup>αε-κλAεβ-βα | |v-λεε |ω-ναε” in IPA “I will go back home.” in English

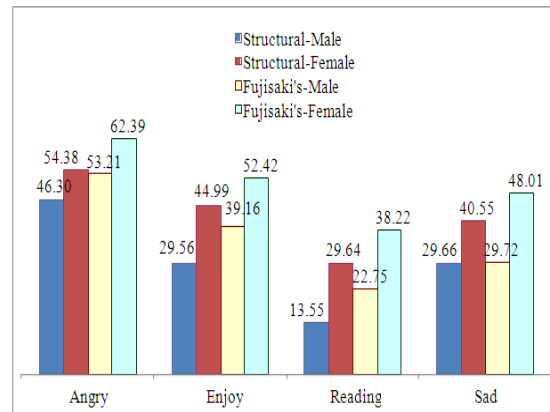


Fig. 9: Averaged RMS error for sentence “χ<sup>η</sup>α&v-ραεκ-κ<sup>η</sup>υν-τΗι | | -συετ-ναφ-λο |κ” in IPA “I love you most in the world.” in English

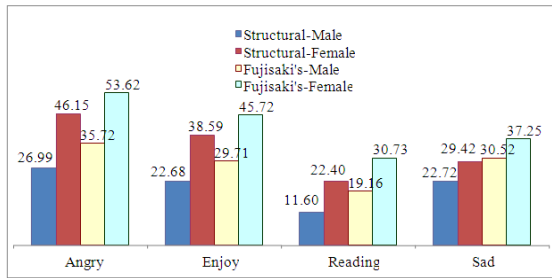


Fig. 10: Averaged RMS error for sentence “ $\rho\alpha\omega\text{-}\mu\alpha\downarrow\phi\text{-}\delta A\downarrow\phi\text{-}\pi\alpha\phi\text{-}\tau H\text{to}\downarrow\omega\text{-}\tau H\alpha\equiv\text{-}\lambda\varepsilon$ ” in IPA “We do not go to the sea.” in English

### DISCUSSION

The experimental results in Fig. 6-10 show that the averaged RMS error of the angry speech is the highest level; meanwhile the averaged RMS error of the reading speech is the lowest level. The averaged RMS errors of the happy and sad speech are in the middle level. It can be obviously seen from all Figures that all 5 sentences have the corresponding results. When considering the differences between genders, we found that the averaged RMS error of female speech is above that of male speech. Last but not least, that the RMS error of the Fujisaki's model is mostly higher than that of the structural model for all speech styles. In other words, the structural model gives the better fit for modeling of the F0 contour of the expressive speech than that of the Fujisaki's model. The results of all sentences confirm this observation (Suhartono, 2011; Souleymane *et al.*, 2009; Geravanchizadeh and Rezaii, 2009; Camminatiello and Lucademo, 2010; Nicknam *et al.*, 2009).

### CONCLUSION

This study proposes a comparison of two successful F0 models. First, the model is based on the Fujisaki's model. Second, the model is based on the structural model. The applied speech database consists of male and female speech and each one contains 4 different speech styles including angry style, sad style, enjoyable style and reading style. Five sentences are used for each speech style and each sentence includes 100 samples. From the experimental results, it has been seen that RMS error of each speech style is different from the others for both models. It also reveals that the RMS error of the Fujisaki's model is higher than that of the structural model for all speech styles. In other words, the structural model gives the better fit for

modeling of the F0 contour of the expressive speech than that of the Fujisaki's model.

### ACKNOWLEDGEMENT

The researcher is grateful to K. Sornthongkam, P. Inthornchai-eur, N. Sangkaew and A. Sricharoenchot for providing the speech databases.

### REFERENCES

- Camminatiello, I. and A. Lucademo, 2010. Estimating multinomial logit model with multicollinear data. *Asian J. Math. Stat.*, 3: 93-101. <http://www.doaj.org/doaj?func=abstract&id=596155>
- Chomphan, S. and T. Kobayashi, 2007. Implementation and evaluation of an HMM-based Thai speech synthesis system. *Proceeding of the 8th Annual Conference of the International Speech Communication Association, (ISCA)*, Aug. 2007, Antwerp, Belgium, pp: 2849-2852. [http://www.isca-speech.org/archive/interspeech\\_2007/i07\\_2849.html](http://www.isca-speech.org/archive/interspeech_2007/i07_2849.html)
- Chomphan, S. and T. Kobayashi, 2008. Tone correctness improvement in speaker dependent HMM-based Thai speech synthesis. *Speech Commun.*, 50: 392-404. DOI: 10.1016/j.specom.2007.12.002
- Chomphan, S. and T. Kobayashi, 2009. Tone correctness improvement in speaker-independent average-voice-based Thai speech synthesis. *Speech Commun.*, 51: 330-343. DOI: 10.1016/j.specom.2008.10.003
- Fujisaki, H. and H. Sudo, 1971. A model for the generation of fundamental frequency contours of Japanese word accent. *J. Acoust. Soc. Jap.*, 57: 445-452. <http://ci.nii.ac.jp/naid/110003107854/en>
- Fujisaki, H., 1983. *Dynamic Characteristics of Voice Fundamental Frequency in Speech and Singing*. In: *The Production of Speech*. Mac Neilage, P.F. (Ed.). Springer, New York, pp: 39-55. [http://www.speech.kth.se/prod/publications/files/qpsr/1981/1981\\_22\\_1\\_001-020.pdf](http://www.speech.kth.se/prod/publications/files/qpsr/1981/1981_22_1_001-020.pdf)
- Geravanchizadeh, M. and T.Y. Rezaii, 2009. Transform domain based multi-channel noise cancellation based on adaptive decorrelation and least mean mixed-norm algorithm. *J. Applied Sci.*, 9: 651-661.
- Hiroya, F. and O. Sumio, 2002. A preliminary study on the modeling of fundamental frequency contours of Thai utterances. *Proceedings of the International Conference on Signal Processing*, Aug. 2002, Beijing, China, pp: 516-519. [http://ieeexplore.ieee.org/xpl/freeabs\\_all.jsp?arnumber=1181106](http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=1181106)



- Li, Y., T. Lee and Y. Qian, 2004. Analysis and modeling of  $F_0$  contours for cantonese text-to-speech. *ACM Trans. Asian Language Inform. Process.*, 3: 169-180. DOI: 10.1145/1037811.1037813
- Ni, J. and K. Hirose, 2006. Quantitative and structural modeling of voice fundamental frequency contours of speech in Mandarin. *Speech Commun.*, 48: 989-1008. DOI: 10.1016/j.specom.2006.01.002
- Nicknam, A., R. Abbasnia, M. Bozorgnasab and Y. Eslamian, 2009. Synthesizing the 2004 Mw 6.2 kojour earthquake using empirical green's function. *Asian J. Sci. Res.*, 2: 119-134. <http://docsdrive.com/pdfs/ansinet/ajsr/2009/119-134.pdf>
- Saito, T. and M. Sakamoto, 2002. Applying a hybrid intonation model to a seamless speech synthesizer. *Proceeding of the International Conference on Spoken Language Processing, Colorado, Sept. 2002, USA.*, pp: 165-168. [http://www.isca-speech.org/archive/icslp\\_2002/i02\\_0165.html](http://www.isca-speech.org/archive/icslp_2002/i02_0165.html)
- Seresangtakul, P. and T. Takara, 2002. Analysis of pitch contour of Thai tone using Fujisaki's model. *Proceeding of the International Conference on Acoustics, Speech and Signal Processing, May 2002, Orlando, USA.*, pp: 505-508. [http://ieeexplore.ieee.org/xpl/freeabs\\_all.jsp?arnumber=1005787](http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=1005787)
- Seresangtakul, P. and T. Takara, 2003. A generative model of fundamental frequency contours for polysyllabic words of Thai tones. *Proceeding of the International Conference on Acoustics, Speech and Signal Processing, Apr. 2003, Hong Kong*, pp: 452-455. [http://ieeexplore.ieee.org/xpl/freeabs\\_all.jsp?arnumber=1198815](http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=1198815)
- Souleymane, K., Z. Tang, S. Dong and Y. Jiang, 2009. Evaluation of some organic pollutants transport into the shallow groundwater and surface water of jiaxing landfill area. *Am. J. Applied Sci.*, 6: 2010-2017. DOI: 10.3844/ajassp.2009.2010.2017.
- Suhartono, 2011. Time series forecasting by using seasonal autoregressive integrated moving average: Subset, multiplicative or additive model. *J. Math. Stat.*, 7: 20-27. DOI: 10.3844/jmssp.2011.20.27.
- Tao J., J. Yu and W. Zhang, 2006. Internal dependence based  $F_0$  model for mandarin TTS system. *Proceedings of the TC-STAR Workshop on Speech-to-Speech Translation, Jun. 2006, Barcelona, Spain*, pp: 171-174. [http://www.elda.org/tcstar-workshop\\_2006/pdfs/tts/tcstar06\\_tao.pdf](http://www.elda.org/tcstar-workshop_2006/pdfs/tts/tcstar06_tao.pdf)
- Tran, D.D., E. Castelli, X. H. Le, J.F. Serignat and V. L. Trinh, 2006. Linear  $F_0$  contour model for Vietnamese tones and Vietnamese syllable synthesis with TD-PSOLA. *Proceeding of the International Symposium on Tonal Aspects of Languages, Apr. 2006, La Rochelle, France*, pp: 137-142 <http://www-mrim.imag.fr/publications/2006/XUA06/TAL2006SubmissionUpdate.pdf>