

<u>Abstract</u>: This paper reviews our ongoing work on generative modeling of speech fundamental frequency (F₀) contours for estimating prosodic features from raw speech data [3]-[6]. Proposed F₀ contour model is formulated by translating "Fujisaki model" into probabilistic model described as discrete-time stochastic process. The motivation behind this formulation is:

- (1) to derive a general parameter estimation framework for Fujisaki model, allowing for introduction of powerful statistical methods
- 2 to construct an automatically trainable version of Fujisaki model so that in future it can be incorporated into statistical-model-based text-to-speech synthesis systems

1. Introduction



- time course of frequency of vocal fold vibration (Manifestation of physical movement of thyroid cartilage)
- contains various types of non-linguistic information
- Speaker's identity, intention, attitude, mood, etc. modeling and analyzing F₀ contours can be potentially useful for many speech applications in which prosodic information plays a significant role. (In speech synthesis, one challenge is to create a natural -sounding pitch contour for the utterance as a whole.)



Fujisaki model" [1]

- F₀ contour typically consists of two components:
- phrase components
- (long term pitch variations over the duration of prosodic units), and
- accent components (short term pitch variations in accented syllables). "Fujisaki model" is a well-founded mathematical model that describes FO contour as the sum of these two components (See below for details).
- A notable feature of this model is that the model parameters are associated with physiologically and linguistically meaningful quantities.
- However, automatic estimation of Fujisaki model parameters from raw F₀ contour has been a difficult task...

Aim of this paper

translate Fujisaki model into a probabilistic model (stochastic process)

- motivation for this is twofold:
- to provide general parameter estimation framework utilizing powerful statistical inference methods, and
- to construct an automatically trainable version of Fujisaki model that can potentially be combined with statistical model for speech synthesis.

2. Original Fujisaki m



 \blacksquare F₀ contour (log scale) =

<u>Phrase component</u> + <u>Accent component</u> + <u>Base value</u>

- contributions associated with forward-backward translation and rotation of thyroid cartilage, respectively
- outputs of different 2nd order critically damped systems



Prior distribution:
$$p(s) = \pi_{s_1} \prod_{k=2}^{K} \pi_{s_{k-1}}$$

