

# An introduction of trajectory model into HMM-based speech synthesis

Heiga Zen Keiich Tokuda Tadashi Kitamura (Nagoya Institute of Technology, Japan)

## Introduction

### Corpus-based speech synthesis system

Unit selection and concatenation (e.g., CHATR)

- High quality (sometimes discontinuous)
- Require large memory
- Difficult to change voice characteristics

Speech synthesis from HMMs (e.g., HTS)

- Vocoded speech (smooth and stable)
- Small footprint
- Easy to change voice characteristics

### Problems in HTS

- Excitation
  - ⇒ Mixed excitation, Sinusoidal, etc.
- Formant structures get smooth
  - ⇒ Post filtering based formant emphasis

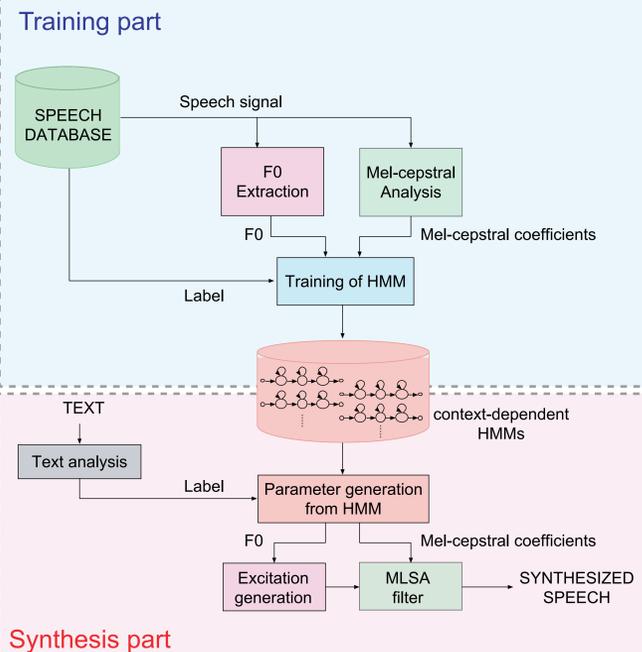
### Trajectory-HMM and its training algorithm

(Tokuda et.al;Eurospeech2003, Zen et.al;ICASSP2004)

- Introduce dynamic feature constraints into HMM
  - ⇒ Reformulate HMM as trajectory model
- Derive its training algorithm based on ML criterion
  - ⇒ Closed-loop training for HTS

⇒ Apply trajectory-HMM training to HTS

## System overview of the HTS



## Dynamic feature constraints

### Speech parameter generation from HMM

⇒ Generate most likely observation sequence

Without dynamic feature constraints

$$o_{\max} = \arg \max_o P(o | q, \lambda)$$

$o$ : observation sequence  
 $q$ : state sequence

⇒  $o_{\max}$  becomes a sequence of mean vectors

Integrate dynamic feature constraints

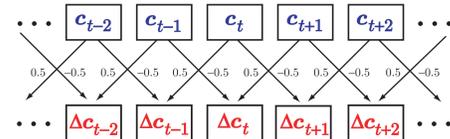
Observation consisted from static and dynamic features

$$o_t = [c_t^T, \Delta c_t^T, \Delta^2 c_t^T]^T$$

$c_t$ : static feature ( $M \times 1$ )  
 $\Delta c_t$ : 1st-order dynamic feature  
 $\Delta^2 c_t$ : 2nd-order dynamic feature

Dynamic feature ⇒ computed from static features

$$\text{Ex.) } \Delta c_t = 0.5c_{t+1} - 0.5c_{t-1}$$



Relationship between static & observation vector sequence

$$o = Wc$$

$W$ : Window matrix ( $3MT \times MT$ )

Under this relationship, we should regard only  $c$  as the random variable

With dynamic feature constraints

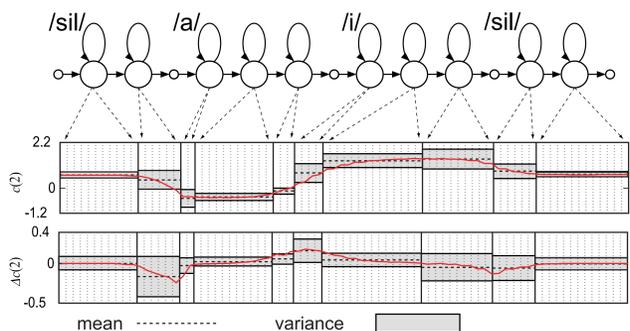
$$c_{\max} = \arg \max_c P(Wc | q, \lambda)$$

By setting  $\frac{\partial}{\partial c} \log P(Wc | q, \lambda) = 0$

$$W^T \Sigma_q^{-1} Wc = W^T \Sigma_q^{-1} \mu_q$$

$\mu_q = [\mu_{q1}, \dots, \mu_{qT}]^T$   
 $\Sigma_q = \text{diag}[\Sigma_{q1}, \dots, \Sigma_{qT}]$

A sequence of speech parameter vector can be determined based on statistics of static and dynamic features.



## Trajectory HMM

### Dynamic feature constraints

- ⇒ Only used in the synthesis stage of HTS
- Inconsistency between training and synthesis

### Trajectory-HMM and its training algorithm

(Tokuda et.al;Eurospeech2003, Zen et.al;ICASSP2004)

- Introduce dynamic feature constraints into HMM
  - ⇒ Reformulate HMM as trajectory model
- Derive its training algorithm based on ML criterion
  - ⇒ Closed-loop training for HTS

⇒ Apply trajectory-HMM training to HTS

### Definition of trajectory-HMM

Output probability for given trajectory-HMM  $\lambda$ :

$$P(c | \lambda) = \sum_{\text{all } q} P(c | q, \lambda) P(q | \lambda)$$

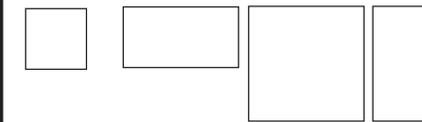
$$P(c | q, \lambda) = \mathcal{N}(c | \bar{c}_q, P_q)$$

$\bar{c}_q$ : Utterance mean vector for given state sequence

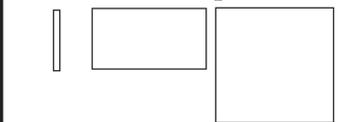
$P_q$ : Utterance covariance matrix for given state sequence

$$\bar{c}_q = P_q r_q$$

$$R_q = W^T \Sigma_q^{-1} W = P_q^{-1}$$



$$r_q = W^T \Sigma_q^{-1} \mu_q$$



⇒  $\bar{c}_q$  is exactly the same to parameter trajectory obtained by speech parameter generation algorithm!

When  $P(c | q, \lambda) = \mathcal{N}(c | \bar{c}_q, P_q)$  is maximized?

⇒ training data  $c$  and generated trajectory  $\bar{c}_q$  (utterance mean vector) is equivalent.

Maximizing model likelihood of the trajectory-HMM = Minimizing the error between training data and generated speech parameter trajectory

⇒ closed-loop training for HTS

## Speech synthesis experiment

Training data	CMU ARCTIC database speaker AWB first 1096 utterances
Test data	Remaining 42 utterances
Sampling rate	16 kHz
Window	25-ms Blackman window
Frame rate	5-ms
Spectral analysis	24-order Mel-cepstral analysis
Dynamic feature	calculated from $\pm 1$ frames
Feature vector	$c(0) \sim c(24)$ , $\log F0$ , and its $\Delta$ , $\Delta\Delta$
Topology	5-state left-to-right HMM with no skip Spectrum: single Gaussian distribution F0: multi-space probability distribution

### Trained models

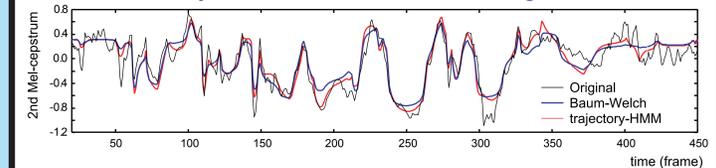
Baum-Welch: Baum-Welch training

Viterbi: Viterbi training

trajectory-HMM: trajectory-HMM training

Prosodies generated from these models were exactly the same

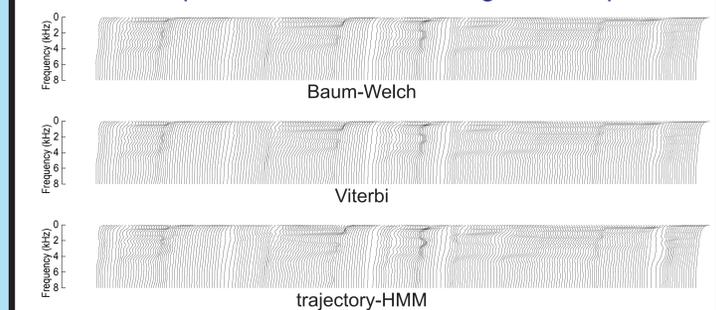
### Generated trajectories and one of training data



Generated trajectory from trajectory-HMM

⇒ Closer to training data than that from Baum-Welch trained HMM

### Generated spectra for a sentence fragment "tropic land"



⇒ Formant structure get clearer slightly

### Subjective listening test

Test type	Paired comparison test
Subjects	8 persons
Test sentences	20 test sentences were chosen at random

