# Statistical Speech Synthesis



Heiga ZEN

Toshiba Research Europe Ltd.

Cambridge Research Laboratory

Speech Synthesis Seminar Series @ CUED, Cambridge, UK

January 11th, 2011

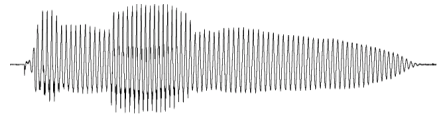# Text-to-speech as a mapping problem

**Text-to-speech synthesis (TTS)**

Text (seq of discrete symbols) → Speech (continuous time series)

**Good morning** →

**Automatic speech recognition (ASR)**

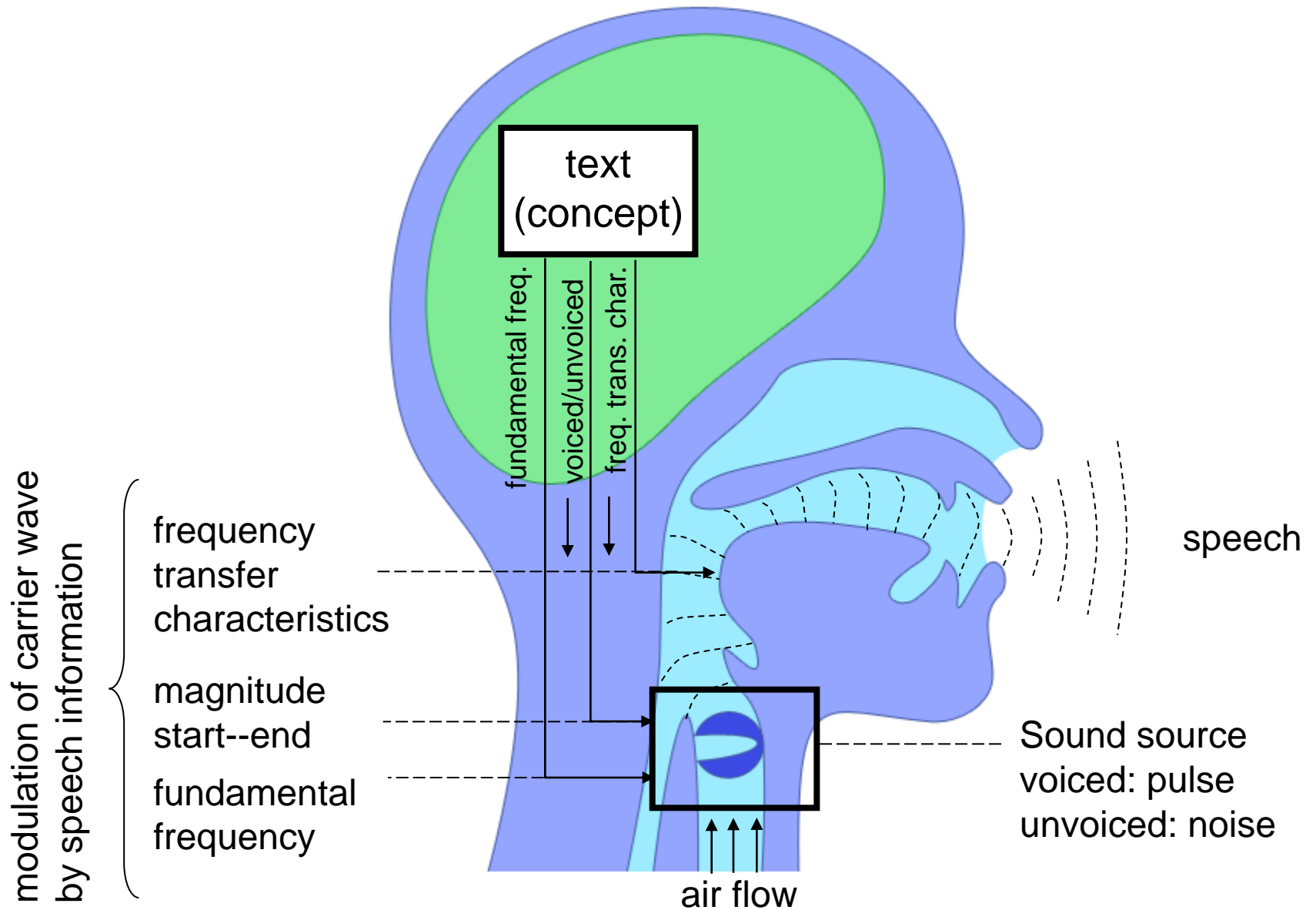Speech (continuous time series) → Text (seq of discrete symbols)

→ **Good morning**

**Machine Translation (MT)**

Text (seq of discrete symbols) → Text (seq of discrete symbols)

**Dobré ráno** → **Good morning**

# Speech production process



text (concept)

fundamental freq.

voiced/unvoiced

freq. trans. char.

modulation of carrier wave by speech information

frequency transfer characteristics

magnitude start--end

fundamental frequency

speech

Sound source
voiced: pulse
unvoiced: noise

air flow

TOSHIBA
Leading Innovation >>>

# Speech synthesis methods (1)

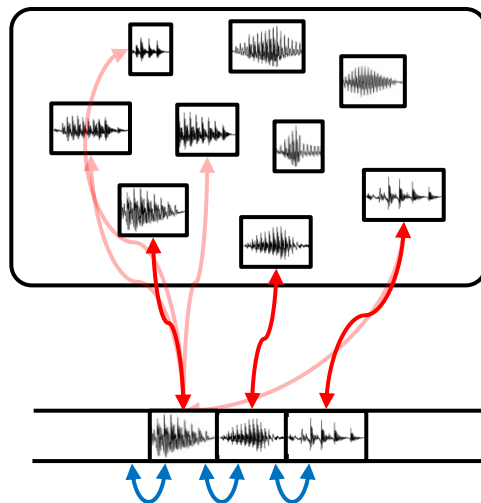## Rule-based, *formant synthesis* (~'90s)



Block diagram of KlattTalk

– Based on parametric representation of speech

– Hand-crafted rules to control phonetic unit

**DECtalk (or KlattTalk / MITTalk) [Klatt;'82]**

# Speech synthesis methods (2)
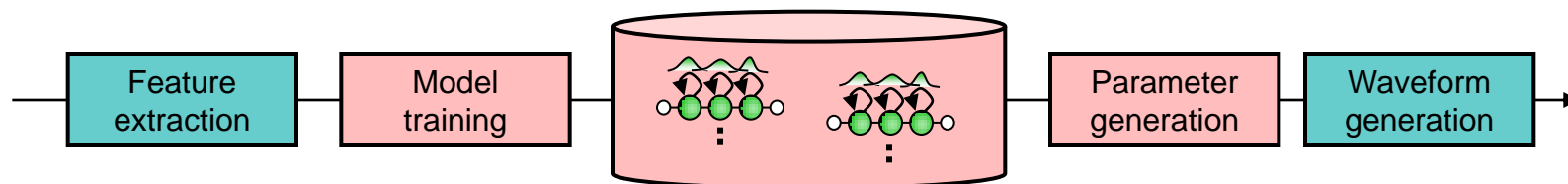
## Corpus-based, *concatenative synthesis* ('90s~)



- – Concatenate small speech units (e.g., phone) from a database
- – Large data + automatic learning → High-quality synthetic voices

**Single inventory; diphone synthesis** [Moullnes;'90]

**Multiple inventory; unit selection synthesis** [Sagisaka;'92, Black;'96]

**TOSHIBA**
Leading Innovation >>>

# Speech synthesis methods (3)

## Corpus-based, *statistical parametric synthesis* (mid '90s~)



- Large data + automatic training

  → Automatic voice building

- Source-filter model + statistical modeling

  → Flexible to change its voice characteristics

**Hidden Markov models (HMMs) as its statistical acoustic model**

→ **HMM-based speech synthesis (HTS)** [Yoshimura;'02]

# Popularity of statistical speech synthesis



# of statistical speech synthesis
related papers in ICASSP

# Aim of this talk

Statistical speech synthesis is getting popular, but…

not many researchers fully understand how it works

Formulate & understand the whole corpus-based speech synthesis process in a unified statistical framework

# Outline

## HMM-based speech synthesis

- Overview

- Implementation of individual components

## Bayesian framework for speech synthesis

- Formulation

- Realizations in HMM-based speech synthesis

- Recent works

## Conclusions

- Summary

- Future research topics

# HMM-based speech synthesis system (HTS)

# HMM-based speech synthesis system (HTS)

# Speech production process



modulation of carrier wave by speech information

frequency transfer characteristics

magnitude start--end

fundamental frequency

speech

Sound source
voiced: pulse
unvoiced: noise

air flow

# Divide speech into frames

Speech is a non-stationary signal

… but can be assumed to be quasi-stationary

→ Divide speech into short-time frames (e.g., 5ms shift, 25ms length)



$$x_1 \quad x_2 \quad x_3 \qquad \qquad x_t \qquad \qquad x_T$$

# Source-filter model

Excitation (source) part    Spectral (filter) part

pulse train
(voiced)                excitation

                        $e(n)$

                                        linear
                                     time-invariant
                                        system              speech

                                         $h(n)$        $x(n) = h(n) * e(n)$

white noise
(unvoiced)

$$x(n) = h(n) * e(n)$$

Fourier transform

$$X(e^{j\omega}) = H(e^{j\omega})E(e^{j\omega})$$

# Spectral (filter) model

**Parametric models speech spectrum**

Autoregressive (AR) model

$$H(z) = c(0) \Bigg/ \left\{ 1 - \sum_{m=1}^{M} c(m) z^{-m} \right\}$$

Exponential (EX) model

$$H(z) = \exp \sum_{m=0}^{M} c(m) z^{-m}$$

ML estimation of spectral model parameters

$$\boldsymbol{c}_t = \arg \max_{\boldsymbol{c}_t} p(\boldsymbol{x}_t \mid \boldsymbol{c}_t)$$

$\boldsymbol{c}_t = [c_t(0), \dots, c_t(M)]^\top$ :
spectral model parameters

$p(\boldsymbol{x}_t \mid \boldsymbol{c}_t)$ : AR model → Linear prediction (LP) [Itakura;'70]

$p(\boldsymbol{x}_t \mid \boldsymbol{c}_t)$ : EX model → ML-based cepstral analysis

# LP analysis (1)

$$\boldsymbol{x}_t = [x_t(1), x_t(2), \ldots, x_t(N)]^\top \quad \text{short-time windowed speech waveform}$$



LP analysis assumes that $\boldsymbol{x}_t$ is a sample from $M$-th order AR process

$$x_t(n) = \sum_{m=1}^{M} c_t(m) x_t(n-m) + \epsilon_t(n) \qquad \text{Linear AR process}$$

$$\boldsymbol{c}_t = [c_t(0), c_t(1), \ldots, c_t(M)]^\top \qquad M\text{-th order LP coefficients}$$

$$\epsilon_t(n) \sim \mathcal{N}(0, c_t(0)) \qquad \text{Gaussian noise}$$

# LP analysis (2)

If we set $\boldsymbol{\Psi}$ as

$$
\boldsymbol{\Psi} =
\begin{bmatrix}
1 & & & & & & 0 \\
-c_t(1) & \ddots & & & & & \\
\vdots & \ddots & \ddots & & & & \\
-c_t(M) & & \ddots & \ddots & & & \\
& \ddots & & \ddots & & \ddots & \\
0 & & & -c_t(M) & \ldots & -c_t(1) & 1
\end{bmatrix}
$$

then

$$
p(\boldsymbol{x}_t \mid \boldsymbol{c}_t) = \mathcal{N}\left(\boldsymbol{x}_t \; ; \; \boldsymbol{0}, c_t(0)\left(\boldsymbol{\Psi}^\top \boldsymbol{\Psi}\right)^{-1}\right)
$$

$$
\hat{\boldsymbol{c}}_t = \arg\max_{\boldsymbol{c}_t} p(\boldsymbol{x}_t \mid \boldsymbol{c}_t) \quad \Rightarrow \text{LP analysis}
$$

# Excitation (source) model



**Excitation model: pulse/noise excitation**

– Voiced (periodic) → pulse trains

– Unvoiced (aperiodic) → white noise

**Excitation model parameters**

– V/UV decision

– V → fundamental frequency (F0): $p_t$

# Speech samples

Natural speech

🔊

Reconstructed speech from extracted parameters (cepstral coefficients & F0 with V/UV decisions)

🔊

**Quality degrades, but main characteristics are preserved**

# HMM-based speech synthesis system (HTS)

# Structure of state-output (observation) vector

$$o_t$$

Spectral parameters
(e.g., cepstrum, LSPs)

$c_t$

$\Delta c_t$ — $\Delta$

Spectrum part

$\Delta^2 c_t$ — $\Delta\Delta$

$p_t$ — log F0 with V/UV

Excitation part

$\Delta p_t$ — $\Delta$

$\Delta^2 p_t$ — $\Delta\Delta$

# Dynamic features

$$\Delta \boldsymbol{c}_t = \frac{\partial \boldsymbol{c}_t}{\partial t} \approx 0.5(\boldsymbol{c}_{t+1} - \boldsymbol{c}_{t-1})$$

$$\Delta^2 \boldsymbol{c}_t = \frac{\partial^2 \boldsymbol{c}_t}{\partial t^2} \approx \boldsymbol{c}_{t+1} - 2\boldsymbol{c}_t + \boldsymbol{c}_{t-1}$$

| $\boldsymbol{c}_{t-1}$ | $\boldsymbol{c}_t$ | $\boldsymbol{c}_{t+1}$ |
|---|---|---|
| $\Delta \boldsymbol{c}_{t-1}$ | $\Delta \boldsymbol{c}_t$ | $\Delta \boldsymbol{c}_{t+1}$ |

| $\boldsymbol{c}_{t-1}$ | $\boldsymbol{c}_t$ | $\boldsymbol{c}_{t+1}$ |
|---|---|---|
| $\Delta^2 \boldsymbol{c}_{t-1}$ | $\Delta^2 \boldsymbol{c}_t$ | $\Delta^2 \boldsymbol{c}_{t+1}$ |

# HMM-based modeling

Label sequence $l$

Sentence HMM $\lambda$

sil     a   i     sil



Observation sequence $o$

State sequence $q$

$o_1$   $o_2$     $\cdots$   $\cdot$   $\cdot$   $\cdots$   $o_T$

1   1   2   3   3     $\cdots$     $N$

# Multi-stream HMM structure

$$b_j(\boldsymbol{o}_t)$$
$$= \prod_{s=1}^{S} \left\{ b_j^s(\boldsymbol{o}_t^s) \right\}^{w_s}$$

Spectrum
(cepstrum or LSP,
& dynamic features)

Excitation
(log F0
& dynamic features)

$\boldsymbol{o}_t$

$\boldsymbol{c}_t$

$\Delta \boldsymbol{c}_t$

$\Delta^2 \boldsymbol{c}_t$

$p_t$

$\Delta p_t$

$\Delta^2 p_t$

$\boldsymbol{o}_t^1$

$\boldsymbol{o}_t^2$

$\boldsymbol{o}_t^3$

$\boldsymbol{o}_t^4$

$b_j(\boldsymbol{o}_t)$

$b_j^1(\boldsymbol{o}_t)$

$b_j^2(\boldsymbol{o}_t)$

$b_j^3(\boldsymbol{o}_t)$

$b_j^4(\boldsymbol{o}_t)$

Stream 1   2   3   4

# Observation of F0



Unable to model by continuous or discrete distribution

# Multi-space probability distribution (MSD)

# Structure of state-output distributions



Stream 1

Spectral params — Single Gaussian

Stream 2,3,4

Log F0

Voiced / Unvoiced — MSD (Gaussian & discrete)

Voiced / Unvoiced — MSD (Gaussian & discrete)

Voiced / Unvoiced — MSD (Gaussian & discrete)

# Training process

data & labels

Compute variance floor

Initialize CI-HMMs by segmental k-means

Reestimate CI-HMMs by EM algorithm

Copy CI-HMMs to CD-HMMs

Reestimate CD-HMMs by EM algorithm

Decision tree-based clustering

Reestimate CD-HMMs by EM algorithm

Untie parameter tying structure

Estimate CD-dur Models from FB stats

Decision tree-based clustering

Estimated dur models

Estimated HMMs

monophone
(context-independent, CI)

fullcontext
(context-dependent, CD)

# HMM-based modeling

# Context-dependent modeling

**Phoneme**

- **current phoneme**
- **{preceding, succeeding} two phonemes**

**Syllable**

- # of phonemes at {preceding, current, succeeding} syllable
- {accent, stress} of {preceding, current, succeeding} syllable
- Position of current syllable in current word
- # of {preceding, succeeding} {accented, stressed} syllable in current phrase
- # of syllables {from previous, to next} {accented, stressed} syllable
- Vowel within current syllable

**Word**

- Part of speech of {preceding, current, succeeding} word
- # of syllables in {preceding, current, succeeding} word
- Position of current word in current phrase
- # of {preceding, succeeding} content words in current phrase
- # of words {from previous, to next} content word

**Phrase**

- # of syllables in {preceding, current, succeeding} phrase

**…..**

Huge # of combinations ⇒ Difficult to have all possible models

# Training process



data & labels

| Compute variance floor |
| Initialize CI-HMMs by segmental k-means |
| Reestimate CI-HMMs by EM algorithm |
| Copy CI-HMMs to CD-HMMs |

monophone
(context-independent, CI)

| Reestimate CD-HMMs by EM algorithm |
| Decision tree-based clustering |
| Reestimate CD-HMMs by EM algorithm |
| Untie parameter tying structure |

fullcontext
(context-dependent, CD)

| Estimate CD-dur Models from FB stats |
| Decision tree-based clustering |

Estimated dur models

Estimated HMMs

# Decision tree-based context clustering [Odell;'95]

# Stream-dependent clustering

**Spectrum & excitation have different context dependency**

**→ Build decision trees separately**

Decision trees
for
mel-cepstrum

Decision trees
for F0

# Training process



data & labels

| Compute variance floor |
| Initialize CI-HMMs by segmental k-means |
| Reestimate CI-HMMs by EM algorithm |
| Copy CI-HMMs to CD-HMMs |

| Reestimate CD-HMMs by EM algorithm |
| Decision tree-based clustering |
| Reestimate CD-HMMs by EM algorithm |
| Untie parameter tying structure |

| Estimate CD-dur Models from FB stats |
| Decision tree-based clustering |

Estimated dur models

Estimated HMMs

monophone
(context-independent, CI)

fullcontext
(context-dependent, CD)

35

# **Estimation of state duration models** [Yoshimura;'98]



$$\chi_{t_0, t_1}(i) \propto \sum_{j \neq i} \alpha_{t_0 - 1}(j) a_{ij} a_{ii}^{t_1 - t_0} \prod_{t=t_0}^{t_1} b_i(\boldsymbol{o}_t)$$

$$\cdot \sum_{k \neq i} a_{ik} b_k(\boldsymbol{o}_{t_1 + 1}) \beta_{t_1 + 1}(k)$$
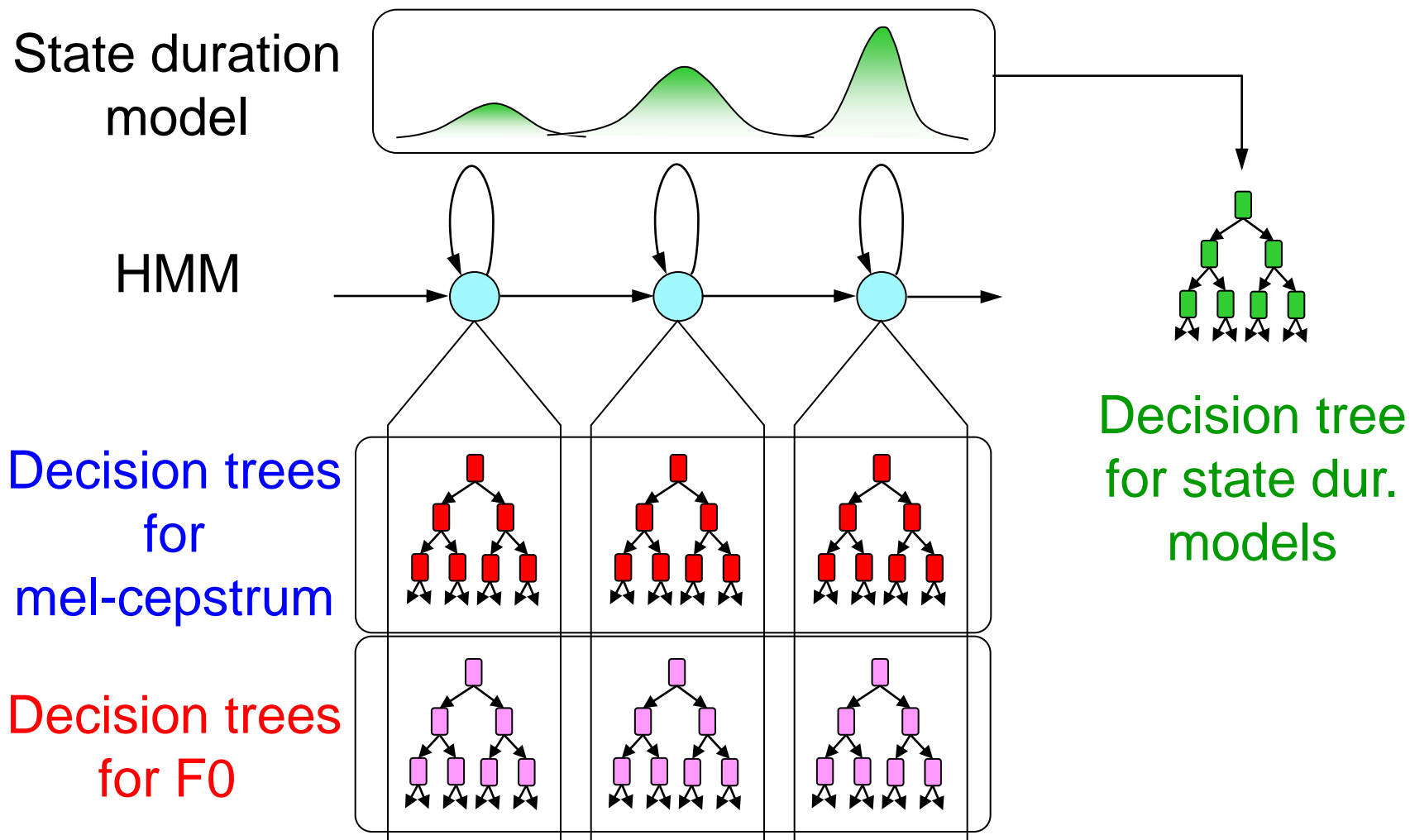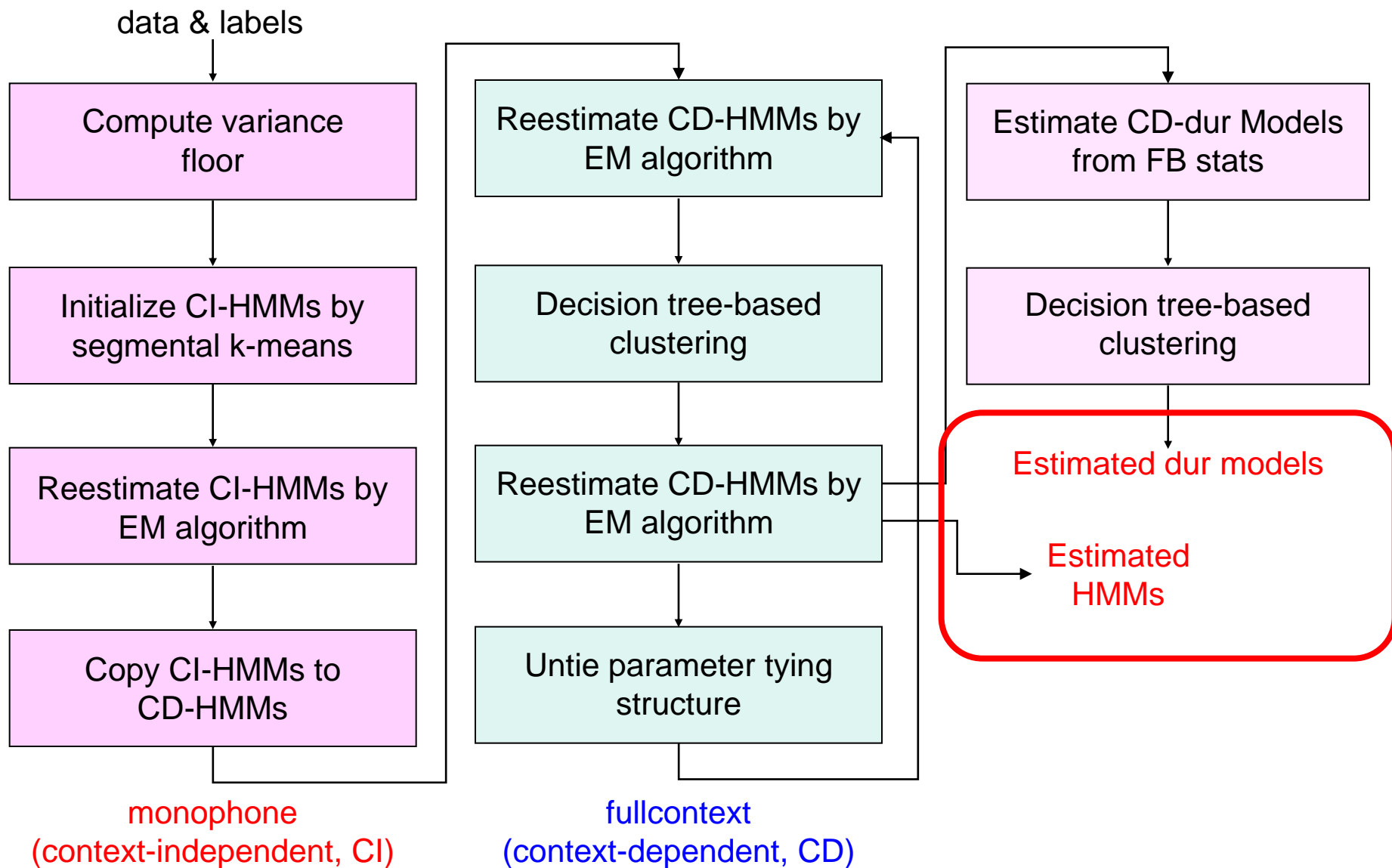
# Stream-dependent clustering



State duration model

HMM

Decision trees for mel-cepstrum

Decision trees for F0

Decision tree for state dur. models

# Training process

data & labels

| monophone (context-independent, CI) | fullcontext (context-dependent, CD) | |
|---|---|---|
| Compute variance floor | Reestimate CD-HMMs by EM algorithm | Estimate CD-dur Models from FB stats |
| Initialize CI-HMMs by segmental k-means | Decision tree-based clustering | Decision tree-based clustering |
| Reestimate CI-HMMs by EM algorithm | Reestimate CD-HMMs by EM algorithm | Estimated dur models |
| Copy CI-HMMs to CD-HMMs | Untie parameter tying structure | Estimated HMMs |

monophone
(context-independent, CI)

fullcontext
(context-dependent, CD)

# HMM-based speech synthesis system (HTS)



39

# Composition of sentence HMM for given text

TEXT

G2P

POS tagging

Text normalization

Pause prediction

Text analysis

context-dependent label sequence $l$

sentence HMM given labels

# Speech parameter generation algorithm

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \boldsymbol{l}, \hat{\lambda})$$

$$= \arg\max_{\boldsymbol{o}} \sum_{\forall \boldsymbol{q}} p(\boldsymbol{o}, \boldsymbol{q} \mid \boldsymbol{l}, \hat{\lambda})$$

$$\approx \arg\max_{\boldsymbol{o}, \boldsymbol{q}} p(\boldsymbol{o}, \boldsymbol{q} \mid \boldsymbol{l}, \hat{\lambda})$$

$$\Downarrow$$

$$\hat{\boldsymbol{q}} = \arg\max_{\boldsymbol{q}} P(\boldsymbol{q} \mid \boldsymbol{l}, \hat{\lambda})$$

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{q}}, \hat{\lambda})$$

41

# Determination of state sequence (1)



Determine state sequence via determining state durations

# Determination of state sequence (2)

$$P(\boldsymbol{q} \mid \boldsymbol{l}, \hat{\lambda}) = \prod_{i=1}^{K} p_i(d_i)$$

$p_i(\cdot)$ : state-duration distribution of $i$-th state

$d_i$ : state duration of $i$-th state

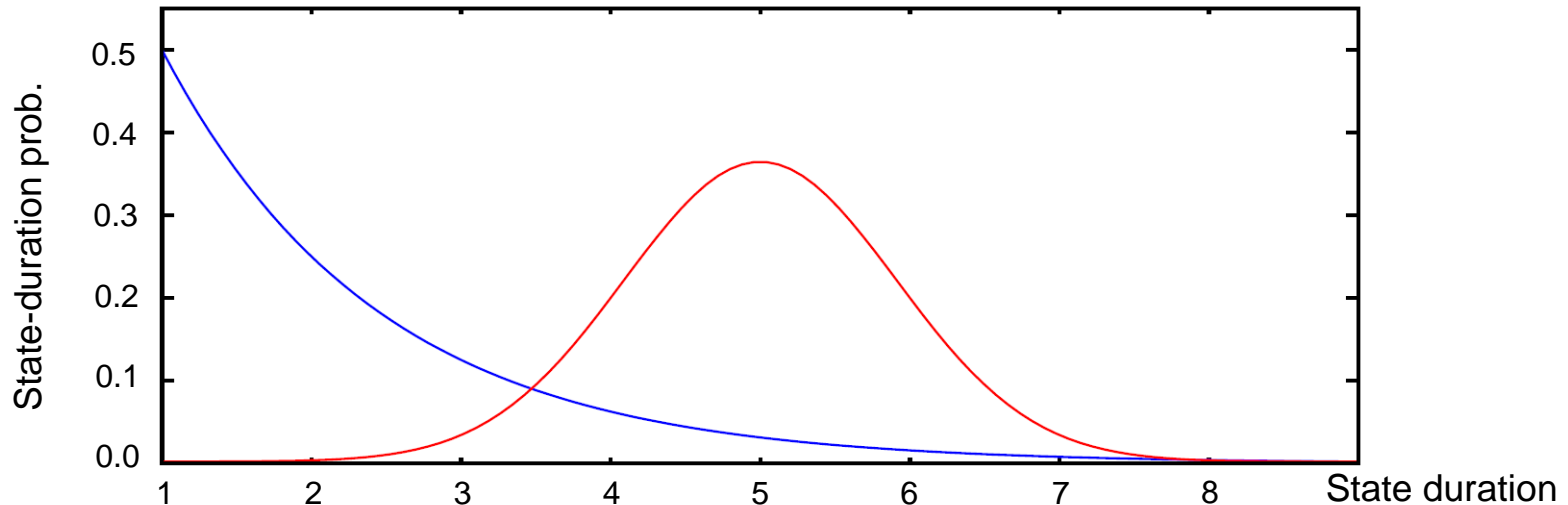$K$ : number of states in a sentence HMM for $w$

# Determination of state sequence (3)

Geometric

$$p_i(d_i) = a_{ii}^{d_i-1}(1 - a_{ii}) \rightarrow \hat{d}_i = \boxed{1}$$

Gaussian

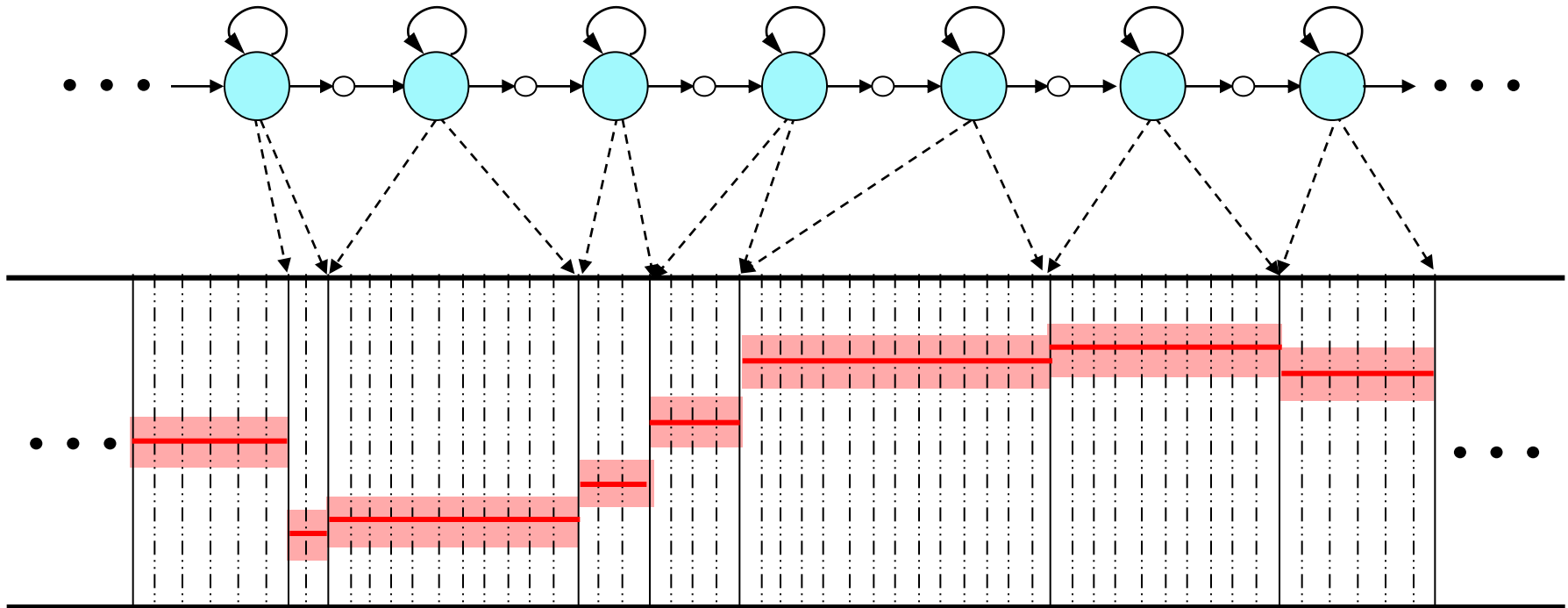$$p_i(d_i) = \mathcal{N}(d_i \; ; \; m_i, \sigma_i^2) \rightarrow \hat{d}_i = \boxed{m_i}$$

# Speech parameter generation algorithm

$$\hat{\boldsymbol{o}} = \arg \max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \boldsymbol{l}, \hat{\lambda})$$

$$= \arg \max_{\boldsymbol{o}} \sum_{\forall \boldsymbol{q}} p(\boldsymbol{o}, \boldsymbol{q} \mid \boldsymbol{l}, \hat{\lambda})$$

$$\approx \arg \max_{\boldsymbol{o}, \boldsymbol{q}} p(\boldsymbol{o}, \boldsymbol{q} \mid \boldsymbol{l}, \hat{\lambda})$$

$\downarrow$

$$\hat{\boldsymbol{q}} = \arg \max_{\boldsymbol{q}} P(\boldsymbol{q} \mid \boldsymbol{l}, \hat{\lambda})$$

$$\hat{\boldsymbol{o}} = \arg \max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{q}}, \hat{\lambda})$$

# Without dynamic features



Mean ——————     Variance ▮

↓

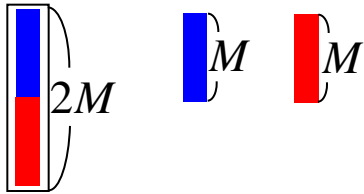$\hat{o}$ → step-wise, mean values

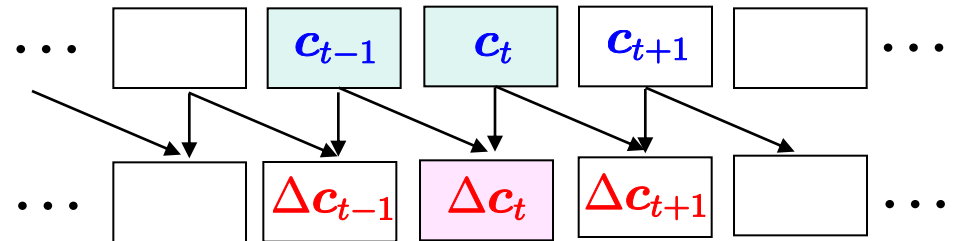# Integration of dynamic features

Speech param. vectors includes both static & dyn. feats.

$$o_t = \left[ c_t^\top, \Delta c_t^\top \right]^\top \qquad \Delta c_t = c_t - c_{t-1}$$



The relationship between $o_t$ & $c_t$ can be arranged as

# Speech parameter generation algorithm

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{q}}, \hat{\lambda})\Big|_{\boldsymbol{o}=\boldsymbol{Wc}}$$

$$\downarrow$$

$$\hat{\boldsymbol{c}} = \arg\max_{\boldsymbol{c}} p(\boldsymbol{Wc} \mid \hat{\boldsymbol{q}}, \hat{\lambda})$$

$$= \arg\max_{\boldsymbol{c}} \mathcal{N}(\boldsymbol{Wc} \, ; \, \boldsymbol{\mu}_{\hat{\boldsymbol{q}}}, \boldsymbol{\Sigma}_{\hat{\boldsymbol{q}}})$$

# Solution

# Generated speech parameter trajectory



Static

Dynamic

Mean ——— Variance ▮ $c$ ———

TOSHIBA
Leading Innovation >>>

# Generated spectra



sil      a      i      sil

w/o dynamic features

w/ dynamic features

# HMM-based speech synthesis system (HTS)

Training part

SPEECH DATABASE

Speech signal

Excitation Parameter extraction

Spectral Parameter Extraction

Excitation parameters

Spectral parameters

Training HMMs

Labels

Context-dependent HMMs & state duration models

TEXT

Text analysis

Labels

Parameter generation from HMMs

Excitation parameters

Spectral parameters

Synthesis part

Excitation generation

Excitation

Synthesis Filter

SYNTHESIZED SPEECH

# Source-filter model

Generated
excitation parameter
(log F0 with V/UV)

Generated
spectral parameter
(cepstrum, LSP)

pulse train

$e(n)$

excitation

white noise

linear
time-invariant
system

$h(n)$

synthesized
speech

$x(n) = h(n) * e(n)$

filtering

# Unvoiced frames & LP spectral coefficients

white noise $\qquad$ Synthesized speech

$$e_t(n) \longrightarrow \boxed{\frac{c_t(0)}{1 - \sum_{m=1}^{M} c_t(m)z^{-m}}} \longrightarrow x_t(n)$$

$$x_t(n) = \sum_{m=1}^{M} c_t(m)x_t(n-m) + e_t(n), \qquad e_t(n) \sim \mathcal{N}(0, c_t(0))$$

**Drive linear filter using white noise**

→ Equivalent to sampling from Gaussian distribution

$$\tilde{\boldsymbol{x}}_t \sim \mathcal{N}\left(\mathbf{0}, c_t(0)\left(\boldsymbol{\Psi}^\top \boldsymbol{\Psi}\right)^{-1}\right)$$

# Speech samples

`w/o dynamic features` 🔊

`w/  dynamic features` 🔊

Use of dynamic features can reduce discontinuity

# Outline

## HMM-based speech synthesis

- Overview
- Implementation of individual components

## Bayesian framework for speech synthesis

- Formulation
- Realizations in HMM-based speech synthesis
- Recent works

## Conclusions

- Summary
- Future research topics

# Statistical framework for speech synthesis (1)

We have a speech database, i.e., a set of texts & corresponding speech waveforms.

Given a text to be synthesized, what is the speech waveform corresponding to the text?

$W$ : set of texts

$X$ : speech waveforms

$w$ : text to be synthesized

$x$ : speech waveform ← unknown

database

Given

**Bayesian framework for prediction**

$$\text{Draw } \tilde{x} \text{ from } p(x \mid w, X, W)$$



$W$ : set of texts

$X$ : speech waveforms

database

Given

$w$ : text to be synthesized

$x$ : speech waveform ⟵ unknown

1. Estimate predictive distribution given variables
2. Draw sample from the distribution

58

# Bayesian framework for speech synthesis (3)

1. Estimating predictive distribution is hard ☹
   → Introduce acoustic model parameters

$$p(\boldsymbol{x} \mid \boldsymbol{w}, \boldsymbol{X}, \boldsymbol{W})$$

⇓ introduce acoustic model $\lambda$

$$= \int p(\boldsymbol{x}, \lambda \mid \boldsymbol{w}, \boldsymbol{W}, \boldsymbol{X})d\lambda = \int p(\boldsymbol{x} \mid \boldsymbol{w}, \lambda)p(\lambda \mid \boldsymbol{W}, \boldsymbol{X})d\lambda$$

$\lambda$ : acoustic model (e.g. HMM )

# Bayesian framework for speech synthesis (4)

2. Using speech waveform directly is difficult ☹
   → Introduce parametric its representation

$$p(\boldsymbol{x} \mid \boldsymbol{w}, \boldsymbol{X}, \boldsymbol{W})$$

$$= \int p(\boldsymbol{x} \mid \boldsymbol{w}, \lambda) p(\lambda \mid \boldsymbol{X}, \boldsymbol{W}) d\lambda$$

$\boldsymbol{x}$     $\boldsymbol{o}$

⇓ introduce parametric representation of speech $\boldsymbol{o}$

$$= \iint p(\boldsymbol{x} \mid \boldsymbol{o}) p(\boldsymbol{o} \mid \boldsymbol{w}, \lambda) p(\lambda \mid \boldsymbol{X}, \boldsymbol{W}) d\lambda d\boldsymbol{o}$$

$\boldsymbol{o}$ : parametric representation of speech waveform $\boldsymbol{x}$
(e.g., cepstrum, LPC, LSP, F0, aperiodicity)

# Bayesian framework for speech synthesis (5)

3. Same texts can have multiple pronunciations, POS, etc. ☹
  → Introduce labels

$$p(\boldsymbol{x} \mid \boldsymbol{w}, \boldsymbol{X}, \boldsymbol{W})$$

$$= \iint p(\boldsymbol{x} \mid \boldsymbol{o})p(\boldsymbol{o} \mid \boldsymbol{w}, \lambda)p(\lambda \mid \boldsymbol{X}, \boldsymbol{W})d\lambda d\boldsymbol{o}$$

⇓ introduce labels derived from texts, $\boldsymbol{l}$ & $\boldsymbol{L}$

$$= \iint \sum_{\forall \boldsymbol{l}} p(\boldsymbol{x} \mid \boldsymbol{o})p(\boldsymbol{o} \mid \boldsymbol{l}, \lambda)P(\boldsymbol{l} \mid \boldsymbol{w})p(\lambda \mid \boldsymbol{X}, \boldsymbol{W})d\lambda d\boldsymbol{o}$$

$\boldsymbol{l}$ : labels derived from text $\boldsymbol{w}$
    (e.g. prons, POS, lexical stress, grammar, pause)

# Bayesian framework for speech synthesis (6)

4. Difficult to perform integral & sum over auxiliary variables ☹
   → Approximated by joint max

$$p(\boldsymbol{x} \mid \boldsymbol{w}, \boldsymbol{X}, \boldsymbol{W})$$

$$= \iint \sum_{\forall \boldsymbol{l}} p(\boldsymbol{x} \mid \boldsymbol{o}) p(\boldsymbol{o} \mid \boldsymbol{l}, \lambda) P(\boldsymbol{l} \mid \boldsymbol{w}) p(\lambda \mid \boldsymbol{X}, \boldsymbol{W}) d\lambda d\boldsymbol{o}$$

⇓ approximate integral & sum by joint max

$$\approx p(\boldsymbol{x} \mid \hat{\boldsymbol{o}}) p(\hat{\boldsymbol{o}} \mid \hat{\boldsymbol{l}}, \hat{\lambda}) P(\hat{\boldsymbol{l}} \mid \boldsymbol{w}) p(\hat{\lambda} \mid \boldsymbol{X}, \boldsymbol{W})$$

where

$$\left\{ \hat{\boldsymbol{o}}, \hat{\boldsymbol{l}}, \hat{\lambda} \right\} = \arg \max_{\boldsymbol{o}, \boldsymbol{l}, \lambda} p(\boldsymbol{x} \mid \boldsymbol{o}) p(\boldsymbol{o} \mid \boldsymbol{l}, \lambda) P(\boldsymbol{l} \mid \boldsymbol{w}) p(\lambda \mid \boldsymbol{X}, \boldsymbol{W})$$

5. Joint maximization is hard ☹
   → Approximated by step-by-step maximizations

$$\left\{\hat{\boldsymbol{o}}, \hat{\boldsymbol{l}}, \hat{\lambda}\right\} = \arg\max_{\boldsymbol{o}, \boldsymbol{l}, \lambda} p(\boldsymbol{x} \mid \boldsymbol{o}) p(\boldsymbol{o} \mid \boldsymbol{l}, \lambda) P(\boldsymbol{l} \mid \boldsymbol{w}) p(\lambda \mid \boldsymbol{X}, \boldsymbol{W})$$

$\Downarrow$ approx joint max by step-by-step max

$$\hat{\lambda} = \arg\max_{\lambda} p(\lambda \mid \boldsymbol{X}, \boldsymbol{W}) \qquad \Leftarrow \text{training}$$

$$\hat{\boldsymbol{l}} = \arg\max_{\boldsymbol{l}} P(\boldsymbol{l} \mid \boldsymbol{w}) \qquad \Leftarrow \text{text analysis}$$

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \hat{\lambda}) \qquad \Leftarrow \text{speech parameter generation}$$

6. Training also requires parametric form of wav & labels ☹
   → Introduce them & approx by step-by-step maximizations

$$\hat{\lambda} = \arg\max_{\lambda} p(\lambda \mid \boldsymbol{X}, \boldsymbol{W})$$

$$\Downarrow$$

$$\hat{\boldsymbol{L}} = \arg\max_{\boldsymbol{L}} P(\boldsymbol{L} \mid \boldsymbol{W}) \qquad \Leftarrow \text{labeling}$$

$$\hat{\boldsymbol{O}} = \arg\max_{\boldsymbol{O}} p(\boldsymbol{X} \mid \boldsymbol{O}) \qquad \Leftarrow \text{feature extraction}$$

$$\hat{\lambda} = \arg\max_{\lambda} p(\hat{\boldsymbol{O}} \mid \hat{\boldsymbol{L}}, \lambda) p(\lambda) \qquad \Leftarrow \text{acoustic model training}$$

$\boldsymbol{O}$ : parametric representation of speech waveforms $\boldsymbol{X}$

$\boldsymbol{L}$ : labels derived from texts $\boldsymbol{W}$

# Bayesian framework for speech synthesis (9)

$$\text{Draw } \tilde{x} \text{ from } p(\boldsymbol{x} \mid \boldsymbol{w}, \boldsymbol{X}, \boldsymbol{W})$$

$$\hat{\boldsymbol{O}} = \arg\max_{\boldsymbol{O}} p(\boldsymbol{X} \mid \boldsymbol{O}) \qquad \Leftarrow \text{feature extraction}$$

$$\hat{\boldsymbol{L}} = \arg\max_{\boldsymbol{L}} P(\boldsymbol{L} \mid \boldsymbol{W}) \qquad \Leftarrow \text{labeling}$$

$$\hat{\lambda} = \arg\max_{\lambda} p(\hat{\boldsymbol{O}} \mid \hat{\boldsymbol{L}}, \lambda) p(\lambda) \qquad \Leftarrow \text{acoustic model training}$$

$$\hat{\boldsymbol{l}} = \arg\max_{\boldsymbol{l}} P(\boldsymbol{l} \mid \boldsymbol{w}) \qquad \Leftarrow \text{text analysis}$$

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \hat{\lambda}) \qquad \Leftarrow \text{speech parameter generation}$$

$$\tilde{\boldsymbol{x}} \text{ from } p(\boldsymbol{x} \mid \hat{\boldsymbol{o}}) \qquad \Leftarrow \text{waveform reconstruction}$$

# HMM-based speech synthesis system (HTS)



**Training part**

SPEECH DATABASE

Speech signal

Excitation extraction

Spectral Extraction

$$\hat{\boldsymbol{O}} = \arg\max_{\boldsymbol{O}} p(\boldsymbol{X} \mid \boldsymbol{O})$$

$$\hat{\boldsymbol{L}} = \arg\max_{\boldsymbol{L}} P(\boldsymbol{L} \mid \boldsymbol{W})$$

Excitation parameters

Spectral parameters

Labels

$$\hat{\lambda} = \arg\max_{\lambda} p(\hat{\boldsymbol{O}} \mid \hat{\boldsymbol{L}}, \lambda) p(\lambda)$$

TEXT

Context-dependent HMMs & state duration models

$$\hat{\boldsymbol{l}} = \arg\max_{\boldsymbol{l}} P(\boldsymbol{l} \mid \boldsymbol{w})$$

Labels

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \hat{\lambda})$$

Excitation parameters

Spectral parameters

Synthesis part

Excitation generation

$$\tilde{\boldsymbol{x}} \sim p(\boldsymbol{x} \mid \hat{\boldsymbol{o}})$$

Synthesis Filter

SYNTHESIZED SPEECH

66

# Problems

## Many approximations

– Integral & sum ≈ max

– Joint max ≈ step-by-step max

    → Poor approximation

## Recent works to relax approximations

– Max → Integral & sum

    ✓ Bayesian acoustic modeling

    ✓ Multiple labels

– Step-wise max → Joint max

    ✓ Statistical vocoding

# Bayesian acoustic modeling (1)

ML-based approach (point estimate of $\lambda$)

$$\hat{\lambda} = \arg\max_{\lambda} p(\hat{\boldsymbol{O}} \mid \hat{\boldsymbol{L}}, \lambda)$$

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \hat{\lambda})$$

Bayesian approach (posterior probability of $\lambda$)

$$\hat{\boldsymbol{o}} = \arg\max_{\boldsymbol{o}} \int p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \hat{\boldsymbol{O}}, \hat{\boldsymbol{L}}) d\lambda$$

$$= \arg\max_{\boldsymbol{o}} \int p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \lambda) p(\lambda \mid \hat{\boldsymbol{O}}, \hat{\boldsymbol{L}}) d\lambda$$

$$= \arg\max_{\boldsymbol{o}} \int p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \lambda) p(\hat{\boldsymbol{O}} \mid \hat{\boldsymbol{L}}, \lambda) p(\lambda) d\lambda$$

**TOSHIBA**
Leading Innovation >>>

# Bayesian acoustic modeling (2)

## Bayesian approach

– Parameters are hidden variables & marginalized out

– Bayesian approach with hidden variables → intractable

→ Variational Bayes [Attias;'99]

$$\log P(\boldsymbol{o}, \hat{\boldsymbol{O}} \mid \hat{\boldsymbol{l}}, \hat{\boldsymbol{L}})$$

$$= \log \sum_{\boldsymbol{q}} \sum_{\boldsymbol{Q}} \int Q(\boldsymbol{q}, \boldsymbol{Q}, \lambda) \frac{P(\boldsymbol{o}, \boldsymbol{q}, \hat{\boldsymbol{O}}, \hat{\boldsymbol{Q}}, \lambda \mid \hat{\boldsymbol{l}}, \hat{\boldsymbol{L}})}{Q(\boldsymbol{q}, \boldsymbol{Q}, \lambda)} d\lambda$$

$$\geq \left\langle \log \frac{P(\boldsymbol{o}, \boldsymbol{q}, \hat{\boldsymbol{O}}, \hat{\boldsymbol{Q}}, \lambda \mid \hat{\boldsymbol{l}}, \hat{\boldsymbol{L}})}{Q(\boldsymbol{q}, \boldsymbol{Q}, \lambda)} d\lambda \right\rangle_{Q(\boldsymbol{q}, \boldsymbol{Q}, \lambda)} \quad \textcolor{red}{\leftarrow \textbf{ Jensen's inequality}}$$

$$= \mathcal{F}$$

# Bayesian acoustic modeling (3)

**Variational Bayesian acoustic modeling for speech synthesis [Nankaku;'03]**

- Fully VB-based speech synthesis
  - ✓ Training posterior distribution of model parameters
  - ✓ Parameter generation from predictive distribution
- Automatic model selection
  - ✓ Bayesian approach provides posterior probability of model structure
- Setting priors
  - ✓ Evidence maximization [Hashimoto;'06]
  - ✓ Cross validation [Hashimoto;'09]
- VB approach works better than ML one when
  - ✓ Data is small
  - ✓ Model is large

# Multiple labels (1)

**Conventional**

$$\hat{\boldsymbol{L}} = \arg \max_{\boldsymbol{L}} P(\boldsymbol{L} \mid \boldsymbol{W}) \qquad\qquad \hat{\lambda} = \arg \max_{\lambda} p(\hat{\boldsymbol{O}} \mid \hat{\boldsymbol{L}}, \lambda) p(\lambda)$$

$$\hat{\boldsymbol{l}} = \arg \max_{\boldsymbol{l}} P(\boldsymbol{l} \mid \boldsymbol{w}) \qquad\qquad \hat{\boldsymbol{o}} = \arg \max_{\boldsymbol{o}} p(\boldsymbol{o} \mid \hat{\boldsymbol{l}}, \hat{\lambda})$$

**Incorporate multiple possible labels**

$$\hat{\lambda} = \arg \max_{\lambda} \sum_{\forall \boldsymbol{L}} p(\hat{\boldsymbol{O}} \mid \boldsymbol{L}, \lambda) P(\boldsymbol{L} \mid \boldsymbol{W}) p(\lambda)$$

$$\hat{\boldsymbol{o}} = \arg \max_{\boldsymbol{o}} \sum_{\forall \boldsymbol{l}} p(\boldsymbol{o} \mid \boldsymbol{l}, \hat{\lambda})$$

**Label sequence is regarded as hidden variable & marginalized**

# Multiple labels (2)

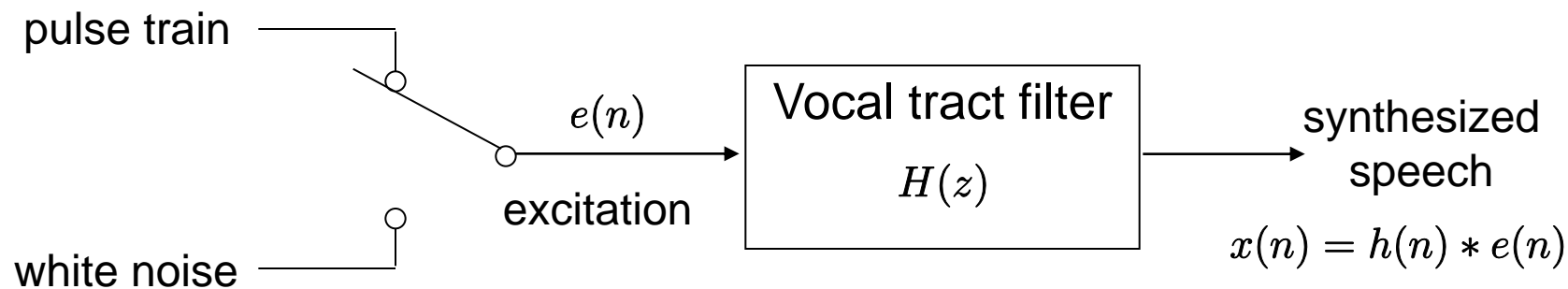**Joint front-end / back-end model training [Oura;'08]**

$$\{\hat{\lambda}, \hat{\Lambda}\} = \arg \max_{\lambda, \Lambda} \sum_{\forall \boldsymbol{L}} p(\hat{\boldsymbol{O}} \mid \boldsymbol{L}, \lambda) P(\boldsymbol{L} \mid \boldsymbol{W}, \Lambda) p(\lambda) p(\Lambda)$$

– Labels = regarded as hidden variable & marginalized

  → Robust against label errors

– Front- & back-end models are trained simultaneously

  → Combine text analysis & acoustic models as a unified model

# Simple pulse/noise vocoding

## Basic pulse/noise vocoder

pulse train

$e(n)$

excitation

Vocal tract filter

$H(z)$
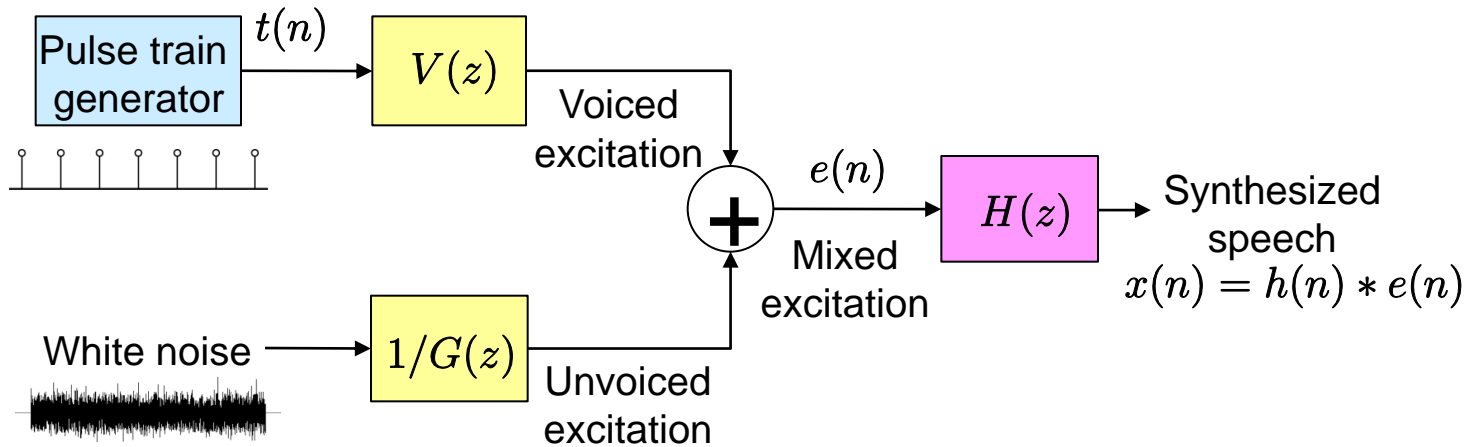
synthesized speech

$x(n) = h(n) * e(n)$

white noise

– Binary switching between voiced & unvoiced excitations

→ Difficult to represent mix of voiced & unvoiced sounds

– Excitations signals of human speech are not pulse or noise

→ Colored voiced/unvoiced excitations

# State-dependent filtering [Maia;'07]

# Waveform-level statistical model (1) [Maia;'10]



$$p(\boldsymbol{x} \mid \boldsymbol{q}, \boldsymbol{c}, \lambda_e) = |\boldsymbol{H_q}|^{-1} \mathcal{N}\left(\boldsymbol{H_q^{-1}} \boldsymbol{x} \; ; \; \boldsymbol{V_q t}, \left(\boldsymbol{G_q^\top G_q}\right)^{-1}\right)$$

$$p(\boldsymbol{x} \mid \boldsymbol{l}, \lambda) = \sum_{\forall \boldsymbol{q}} \int p(\boldsymbol{x} \mid \boldsymbol{q}, \boldsymbol{c}, \lambda_e) p(\boldsymbol{c} \mid \boldsymbol{q}, \lambda_c) p(\boldsymbol{q} \mid \boldsymbol{l}, \lambda_c) d\boldsymbol{c}$$

$$p(\boldsymbol{x} \mid \boldsymbol{w}, \lambda) = \sum_{\forall \boldsymbol{l}} p(\boldsymbol{x} \mid \boldsymbol{l}, \lambda) p(\boldsymbol{l} \mid \boldsymbol{w}) \quad \Leftarrow \text{ waveform-level statistical model}$$

$\boldsymbol{H_q}, \boldsymbol{V_q}, \boldsymbol{G_q}$    matrices representing impulse responses of $H(z), V(z)$, & $G(z)$

$\lambda = \{\lambda_e, \lambda_c\}$    set of acoustic ($\lambda_c$) & excitation ($\lambda_e$) model parameters

75

# Waveform-level statistical model (2) [Maia;'10]

Integral & sum are intractable ☹

→ Approx integral & sum by joint max

$$p(\boldsymbol{x} \mid \boldsymbol{l}, \lambda) = \sum_{\forall \boldsymbol{q}} \int p(\boldsymbol{x} \mid \boldsymbol{q}, \boldsymbol{c}, \lambda) p(\boldsymbol{c} \mid \boldsymbol{q}, \lambda) p(\boldsymbol{q} \mid \boldsymbol{l}, \lambda) d\boldsymbol{c}$$

$$\approx p(\boldsymbol{x} \mid \hat{\boldsymbol{q}}, \hat{\boldsymbol{c}}, \lambda) p(\hat{\boldsymbol{c}} \mid \hat{\boldsymbol{q}}, \lambda) p(\hat{\boldsymbol{q}} \mid \boldsymbol{l}, \lambda) = p(\boldsymbol{x}, \hat{\boldsymbol{q}}, \hat{\boldsymbol{c}} \mid \boldsymbol{l}, \lambda)$$

iteratively optimize $\lambda$ & $C$

$$\hat{\boldsymbol{C}} = \arg\max_{\boldsymbol{C}} p(\boldsymbol{X} \mid \hat{\boldsymbol{Q}}, \boldsymbol{C}, \hat{\lambda}) p(\boldsymbol{C} \mid \hat{\boldsymbol{Q}}, \hat{\lambda}) p(\hat{\boldsymbol{Q}} \mid \hat{\boldsymbol{L}}, \hat{\lambda}) \quad \Leftarrow \text{estimate } C \text{ given } \hat{\lambda}$$

$$\hat{\lambda} = \arg\max_{\lambda} p(\boldsymbol{X} \mid \hat{\boldsymbol{Q}}, \hat{\boldsymbol{C}}, \lambda) p(\hat{\boldsymbol{C}} \mid \hat{\boldsymbol{Q}}, \lambda) p(\hat{\boldsymbol{Q}} \mid \hat{\boldsymbol{L}}, \lambda) \quad \Leftarrow \text{estimate } \lambda \text{ given } \hat{C}$$

Conventional → step-by-step maximization
Proposed → iterative joint maximization

# Outline

## HMM-based speech synthesis

- – Overview
- – Implementation of individual components

## Bayesian framework for speech synthesis

- – Formulation
- – Realizations in HMM-based speech synthesis
- – Recent works

## Conclusions

- – Summary
- – Future research topics

# Summary

**HMM-based speech synthesis**

- Statistical parametric speech synthesis approach
- Source-filter representation of speech + statistical acoustic modeling
- Getting popular

**Bayesian framework for speech synthesis**

- Formulation
- Decomposition to sub-problems
- Correspondence between sub-problems & modules in HMM-based speech synthesis system
- Recent works to relax approximations

# Drawbacks of HMM-based speech synthesis

**Quality of synthesized speech**

- Buzzy
- Flat
- Muffled

**Three major factors degrade the quality**

- Poor vocoding

  → how to parameterize speech?

- Inaccurate acoustic modeling

  → how to model extracted speech parameter trajectories?

- Over-smoothing

  → how to recover generated speech parameter trajectories?

**Still need a lot of works to improve the quality**

# Future challenging topics in speech synthesis

**Keynote speech by Simon King in ISCA SSW7 last year**

Speech synthesis is easy, if ...

- voice is built offline & carefully checked for errors

- speech is recorded in clean conditions

- word transcriptions are correct

- accurate phonetic labels are available or can be obtained

- speech is in the required language & speaking style

- speech is from a suitable speaker

- a native speaker is available, preferably a linguist

**Speech synthesis is not easy if we don't have right data**

# Future challenging topics in speech synthesis

**Non-professional speakers**

- AVM + adaptation (CSTR)

**Too little speech data**

- VTLN-based rapid speaker adaptation (Titech, IDIAP)

**Noisy recordings**

- Spectral subtraction & AVM + adaptation (CSTR)

**No labels**

- Un- / Semi-supervised voice building (CSTR, NICT, CMU, Toshiba)

**Insufficient knowledge of the language or accent**

- Letter (grapheme)-based synthesis (CSTR)
- No prosodic contexts (CSTR, Titech)

**Wrong language**

- Cross-lingual speaker adaptation (MSRA, EMIME)
- Speaker & language adaptive training (Toshiba)

# Thanks!

TOSHIBA
Leading Innovation >>>