

Speech Synthesis as A Machine Learning Problem

Keiichi Tokuda
Nagoya Institute of Technology

Introduction

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

- Hand-crafting each phonetic units by rules

Corpus-based, *concatenative synthesis* ('90s~)

- Concatenate speech units (waveform) from a database
 - Single inventory: diphone synthesis
 - Multiple inventory: unit selection synthesis

Corpus-based, *statistical parametric synthesis*

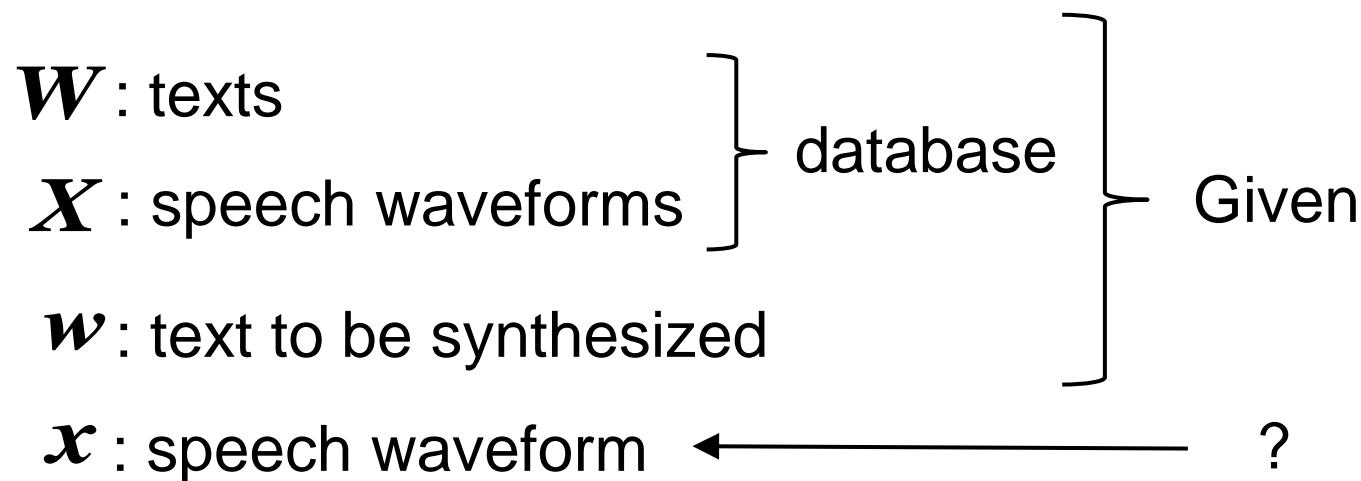
- Source-filter model + statistical acoustic model
 - Hidden Markov model: **HMM-based synthesis**

How we can formulate and understand the whole corpus-based speech synthesis process in a unified statistical framework?

Problem of speech synthesis (1/2)

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

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



Problem of speech synthesis (2/2)

Problem of statistical parametric speech synthesis:

$$\begin{aligned}\hat{\boldsymbol{o}} &= \arg \max_{\boldsymbol{o}} P(\boldsymbol{o} | \boldsymbol{l}, \boldsymbol{O}, \boldsymbol{L}) \\ &= \arg \max_{\boldsymbol{o}} \int \underline{P(\boldsymbol{o} | \boldsymbol{l}, \boldsymbol{\lambda})} \underline{P(\boldsymbol{\lambda} | \boldsymbol{O}, \boldsymbol{L})} d\boldsymbol{\lambda}\end{aligned}$$

\boldsymbol{L} : Label sequence for training data
(pronunciation, stress, POS, pause position, etc.)

\boldsymbol{O} : Training data (speech parameter sequence)

\boldsymbol{l} : Label sequence for synthesis data

\boldsymbol{o} : Synthesis data (speech parameter sequence)

$\boldsymbol{\lambda}$: Model parameters (HMM parameters)

Approximation

Maximum A Posterior (MAP) & ML estimation

$$\hat{\lambda}_{\text{MAP}} = \arg \max_{\lambda} P(\lambda | O, L)$$

MAP estimation

$$= \arg \max_{\lambda} P(O | L, \lambda)P(\lambda)$$

$$\hat{\lambda}_{\text{ML}} = \arg \max_{\lambda} P(O | L, \lambda)$$

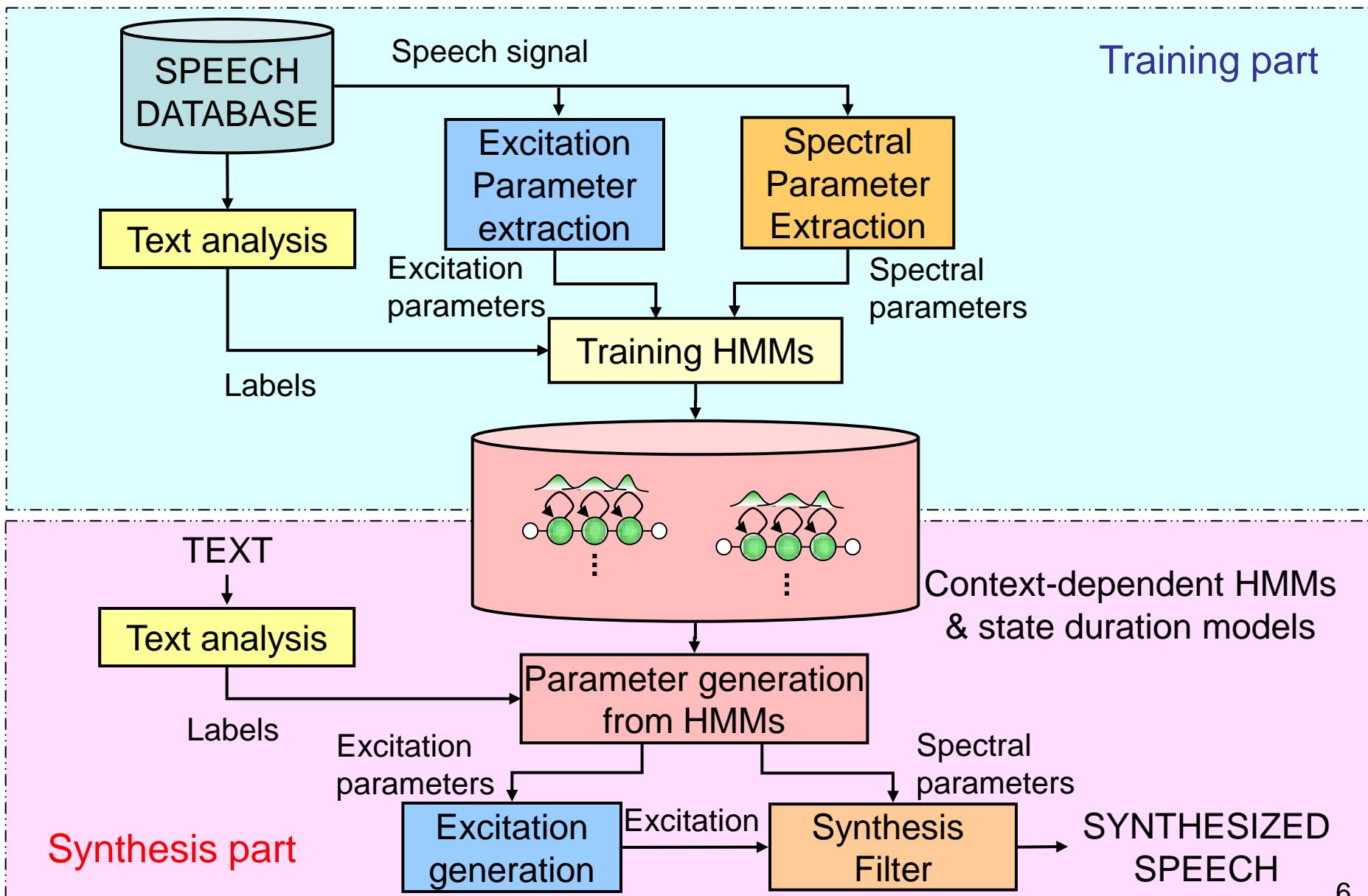
ML estimation

Speech parameter generation using estimated parameters

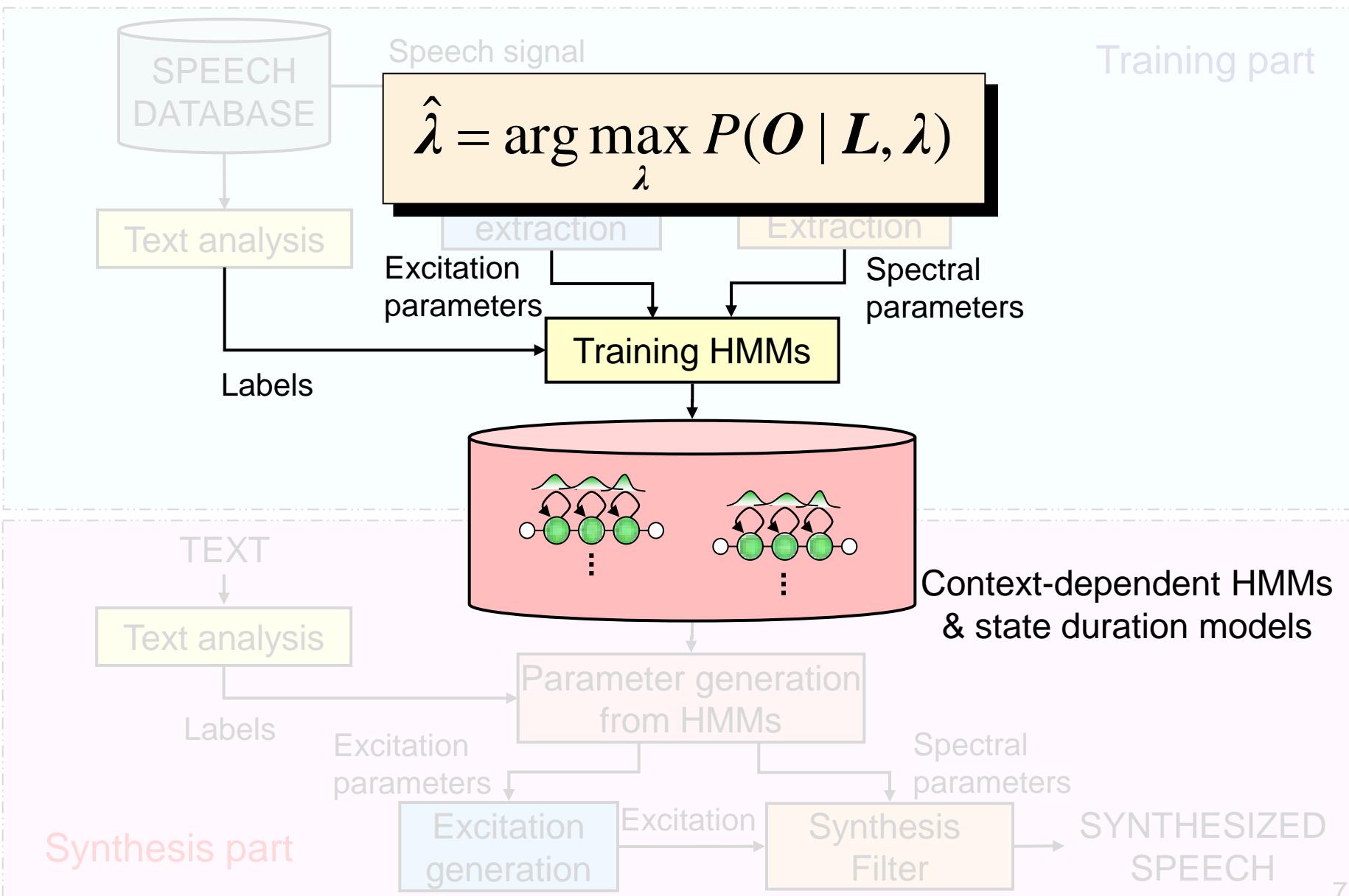
$$\hat{o} = \arg \max_o P(o | l, \hat{\lambda})$$

Speech parameter generation

HMM-based speech synthesis system



HMM-based speech synthesis system

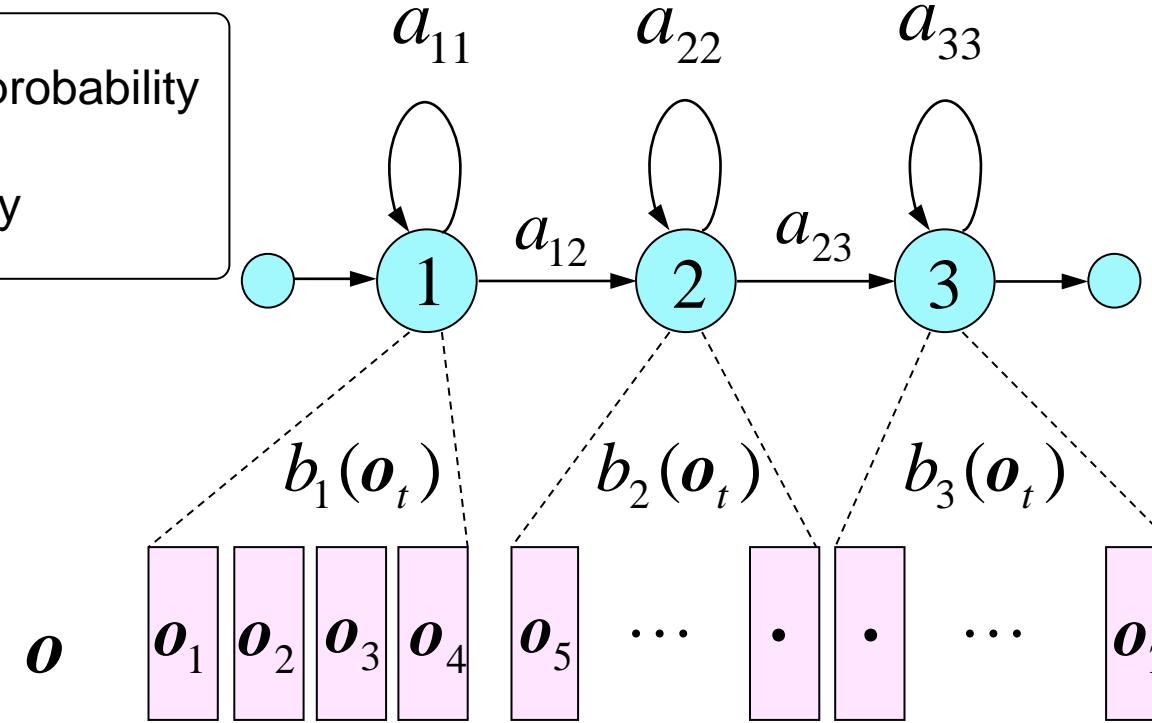


Hidden Markov model (HMM)

a_{ij} : state transition probability

$b_q(o_t)$: output probability

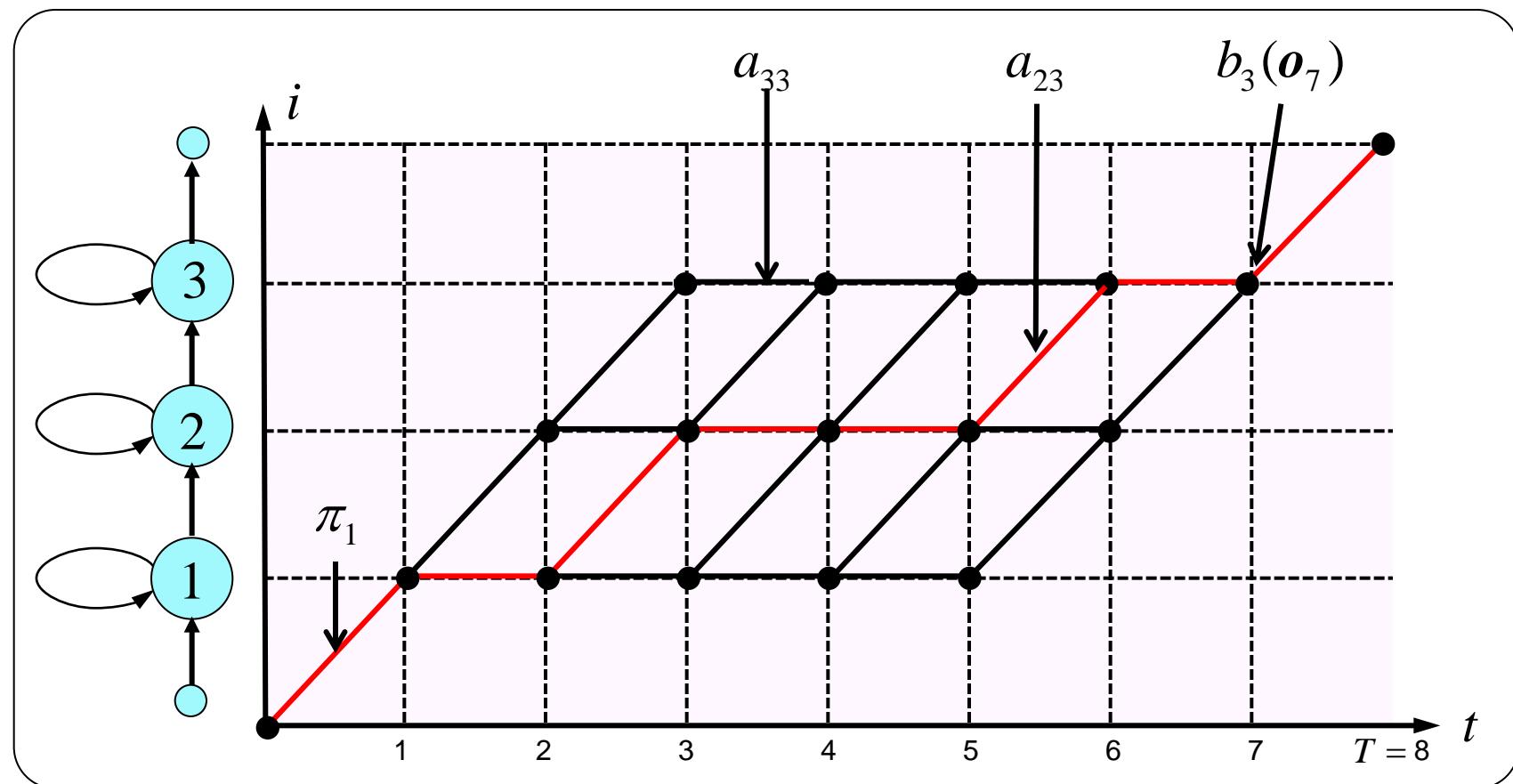
Observation sequence



State sequence

q 1 1 1 1 2 ... 2 3 ... 3

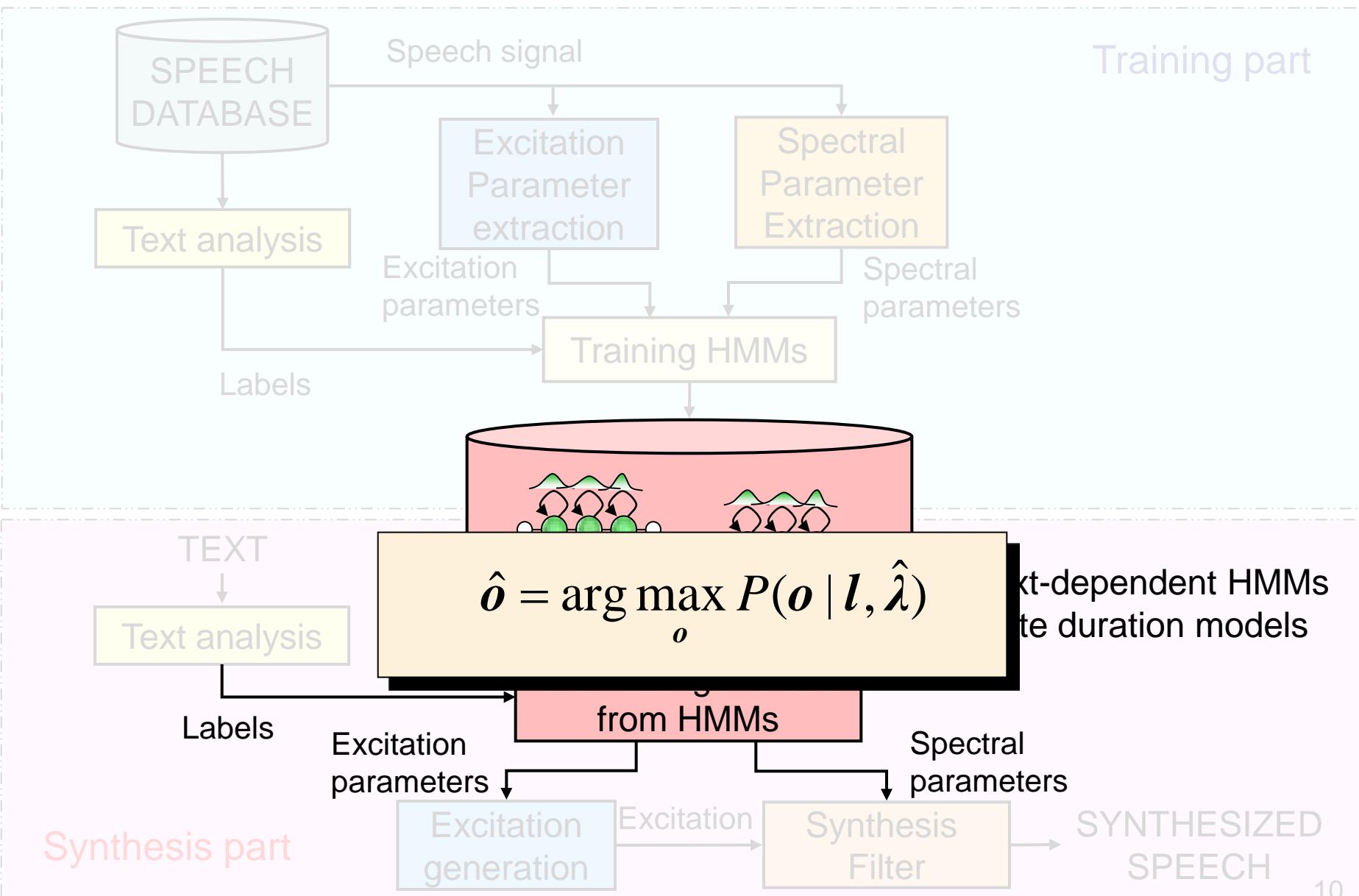
Output probability of HMM



$$P(\mathbf{o} | \lambda) = \sum_{\mathbf{q}} P(\mathbf{o}, \mathbf{q} | \lambda) = \sum_{\mathbf{q}} \prod_{t=1}^T a_{q_{t-1} q_t} b_{q_t} (\mathbf{o}_t)$$

Model parameters can be estimated by an EM algorithm

HMM-based speech synthesis system



Speech parameter generation algorithm

For given HMM λ , determine a speech parameter vector sequence $\mathbf{o} = [\mathbf{o}_1^\top, \mathbf{o}_2^\top, \dots, \mathbf{o}_T^\top]^\top$ which maximizes

$$\begin{aligned} P(\mathbf{o} | \lambda) &= \sum_{\mathbf{q}} P(\mathbf{o} | \mathbf{q}, \lambda) P(\mathbf{q} | \lambda) \\ &\approx \max_{\mathbf{q}} P(\mathbf{o} | \mathbf{q}, \lambda) P(\mathbf{q} | \lambda) \end{aligned}$$



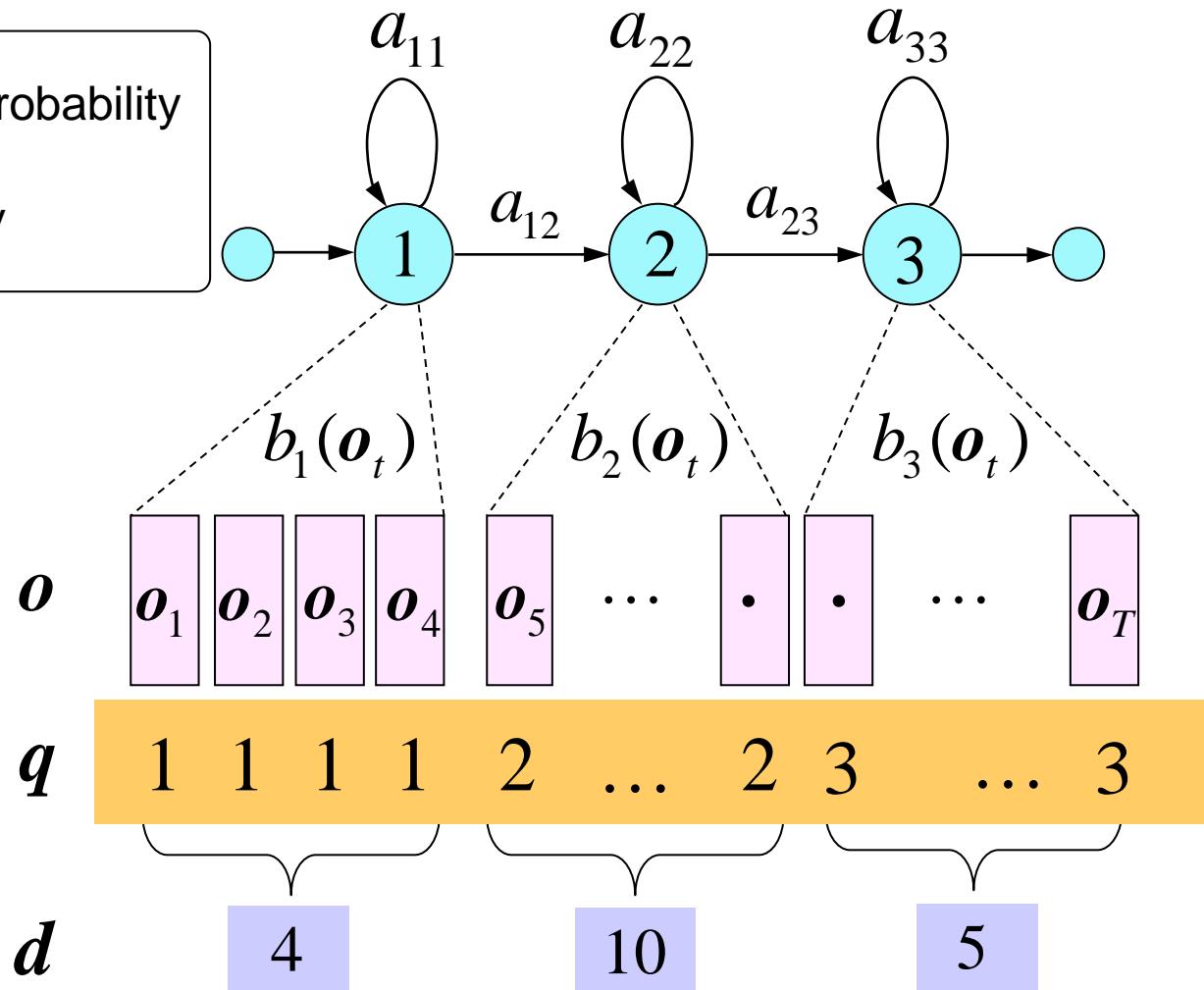
$$\hat{\mathbf{q}} = \arg \max_{\mathbf{q}} p(\mathbf{q} | \mathbf{l}, \lambda)$$

$$\hat{\mathbf{o}} = \arg \max_{\mathbf{o}} p(\mathbf{o} | \hat{\mathbf{q}}, \lambda)$$

Determination of state sequence

a_{ij} : state transition probability

$b_q(o_t)$: output probability



The state sequence can be determined by state durations

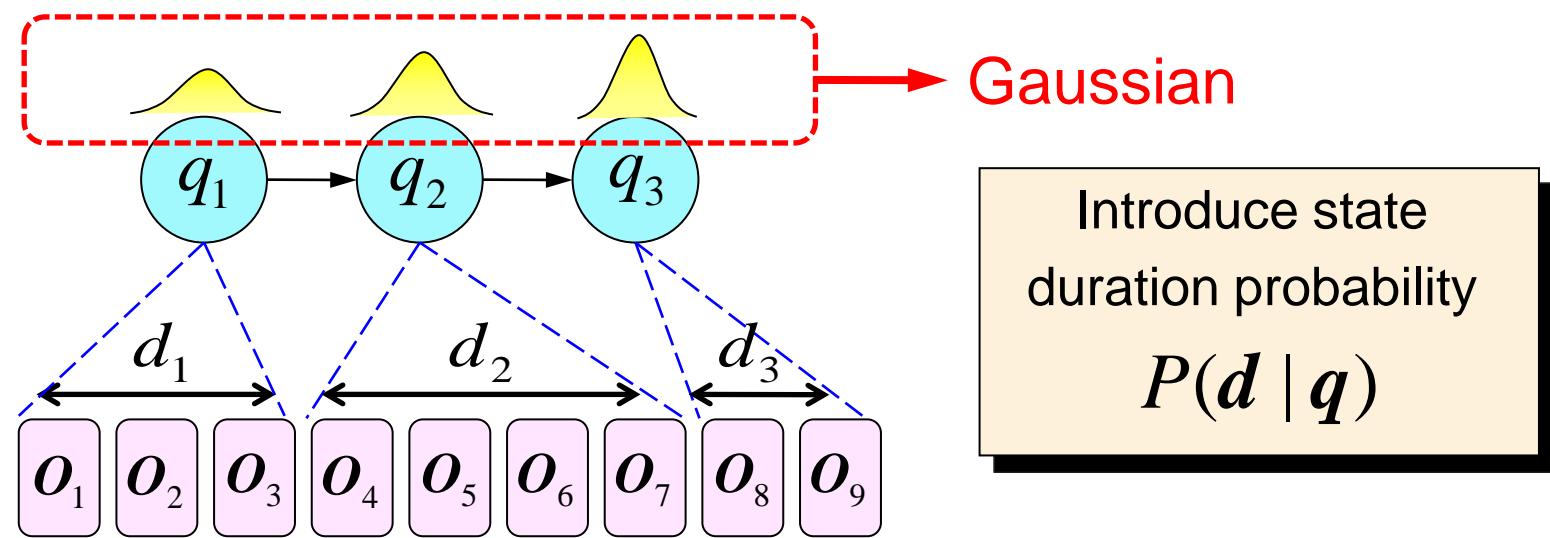
Hidden Semi Markov Model

HMM (Hidden Markov Model)

- State duration prob. depends only on transition prob.
- State duration probability exponentially decreases

HSMM (Hidden Semi Markov Model)

- HMM + explicit duration model \Rightarrow HSMM



State durations are given by means of Gaussians.

Speech parameter generation algorithm

For given HMM λ , determine a speech parameter vector Sequence $\mathbf{o} = [\mathbf{o}_1^\top, \mathbf{o}_2^\top, \dots, \mathbf{o}_T^\top]^\top$ which maximizes

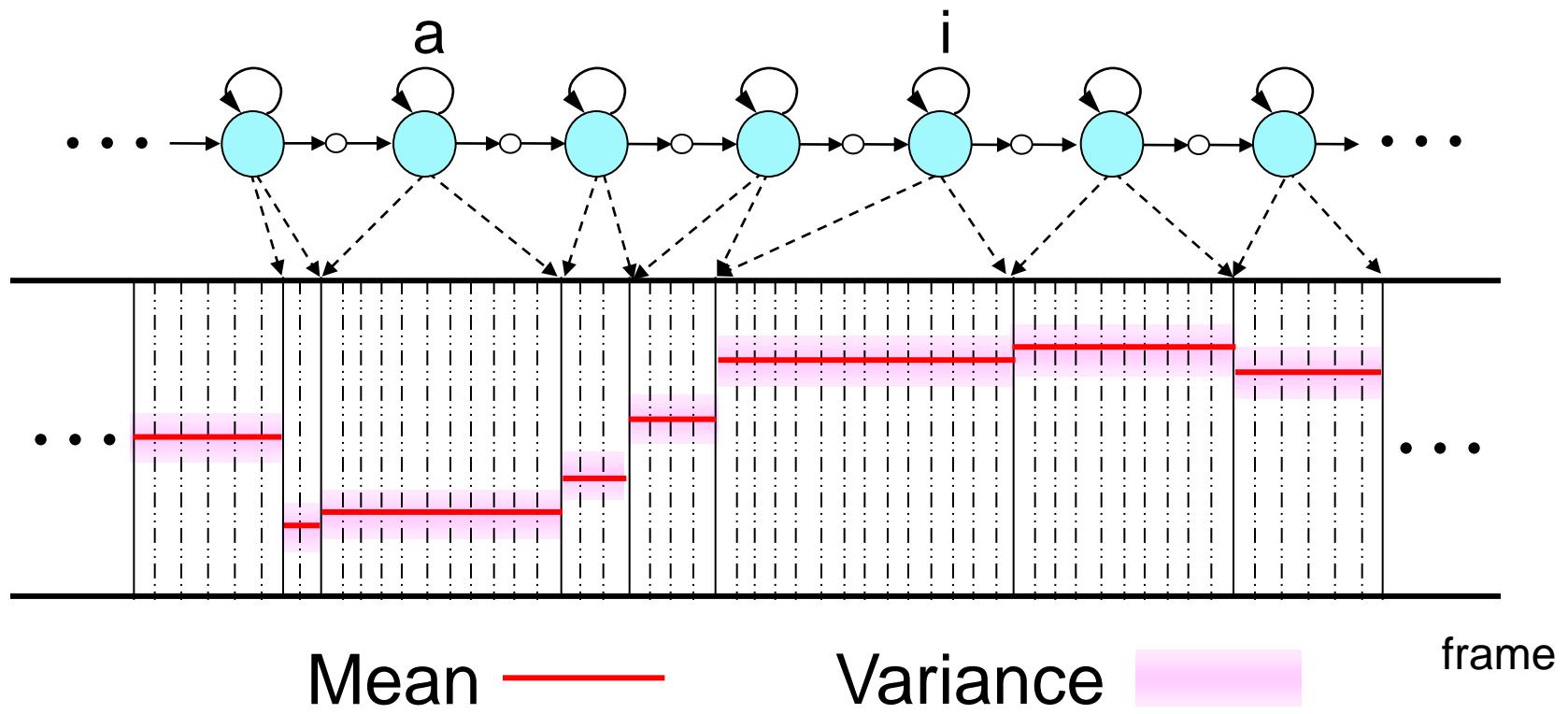
$$\begin{aligned} P(\mathbf{o} | \lambda) &= \sum_{\mathbf{q}} P(\mathbf{o} | \mathbf{q}, \lambda) P(\mathbf{q} | \lambda) \\ &\approx \max_{\mathbf{q}} P(\mathbf{o} | \mathbf{q}, \lambda) P(\mathbf{q} | \lambda) \end{aligned}$$



$$\hat{\mathbf{q}} = \arg \max_{\mathbf{q}} p(\mathbf{q} | \mathbf{l}, \lambda)$$

$$\hat{\mathbf{o}} = \arg \max_{\mathbf{o}} p(\mathbf{o} | \hat{\mathbf{q}}, \lambda)$$

Generated feature sequence

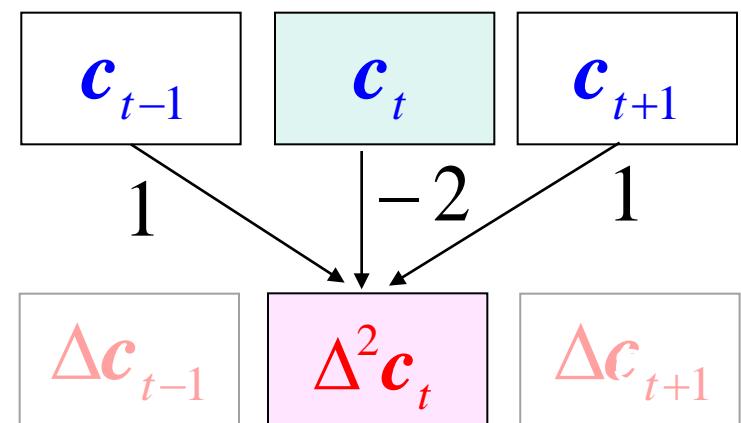
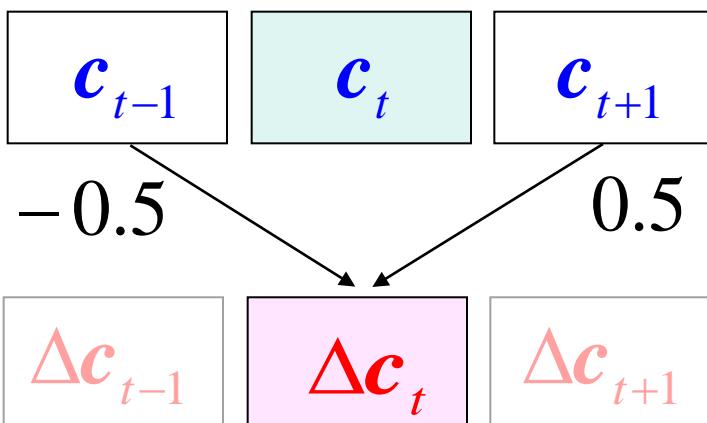


\hat{o} becomes a sequence of mean vectors
⇒ discontinuous outputs between states

Dynamic features

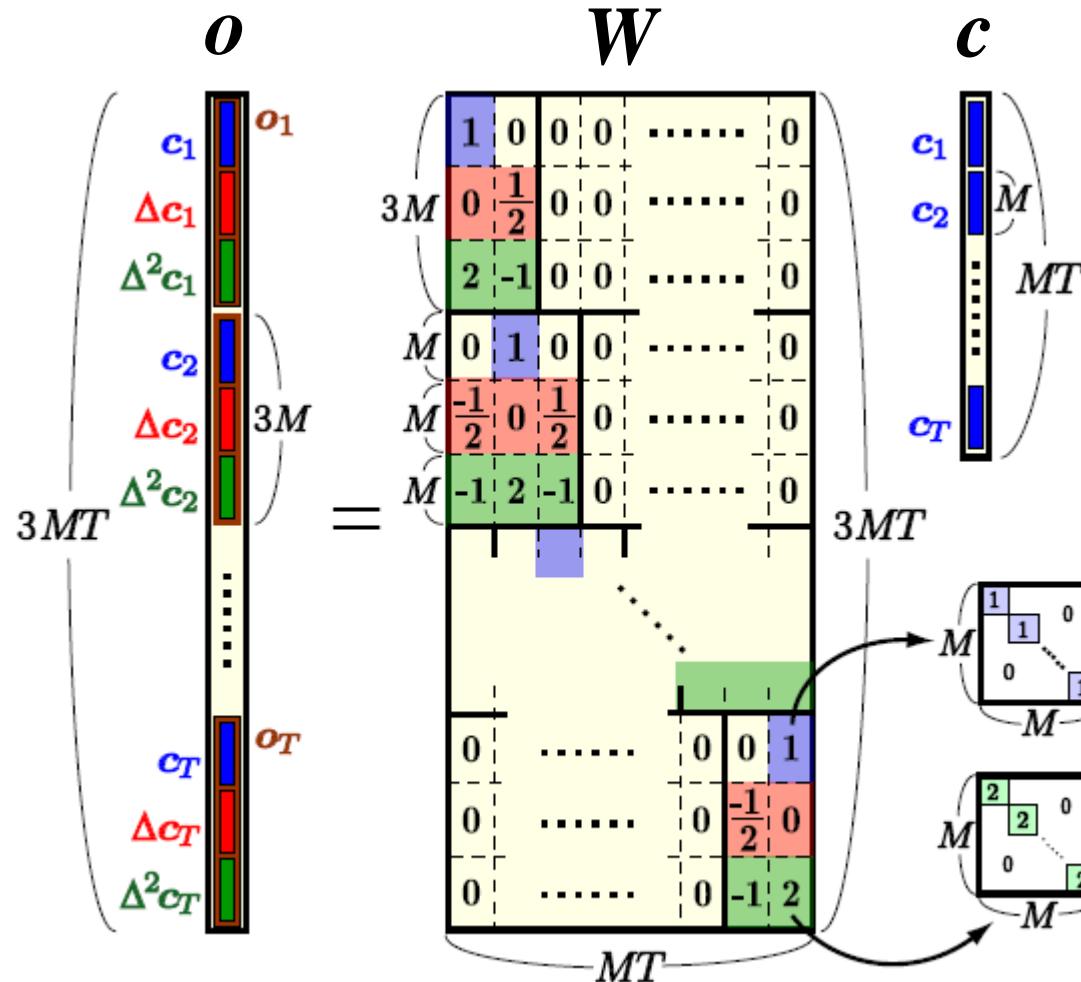
$$\Delta \mathbf{c}_t = \frac{\partial \mathbf{c}_t}{\partial t} \approx 0.5(\mathbf{c}_{t+1} - \mathbf{c}_{t-1})$$

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



Integration of dynamic features

Relationship between speech parm. vec. & static feat



Solution for the problem (1/2)

By setting

$$\frac{\partial \log P(\mathbf{W}\mathbf{c} | \hat{\mathbf{q}}, \lambda)}{\partial \mathbf{c}} = \mathbf{0}$$

we obtain

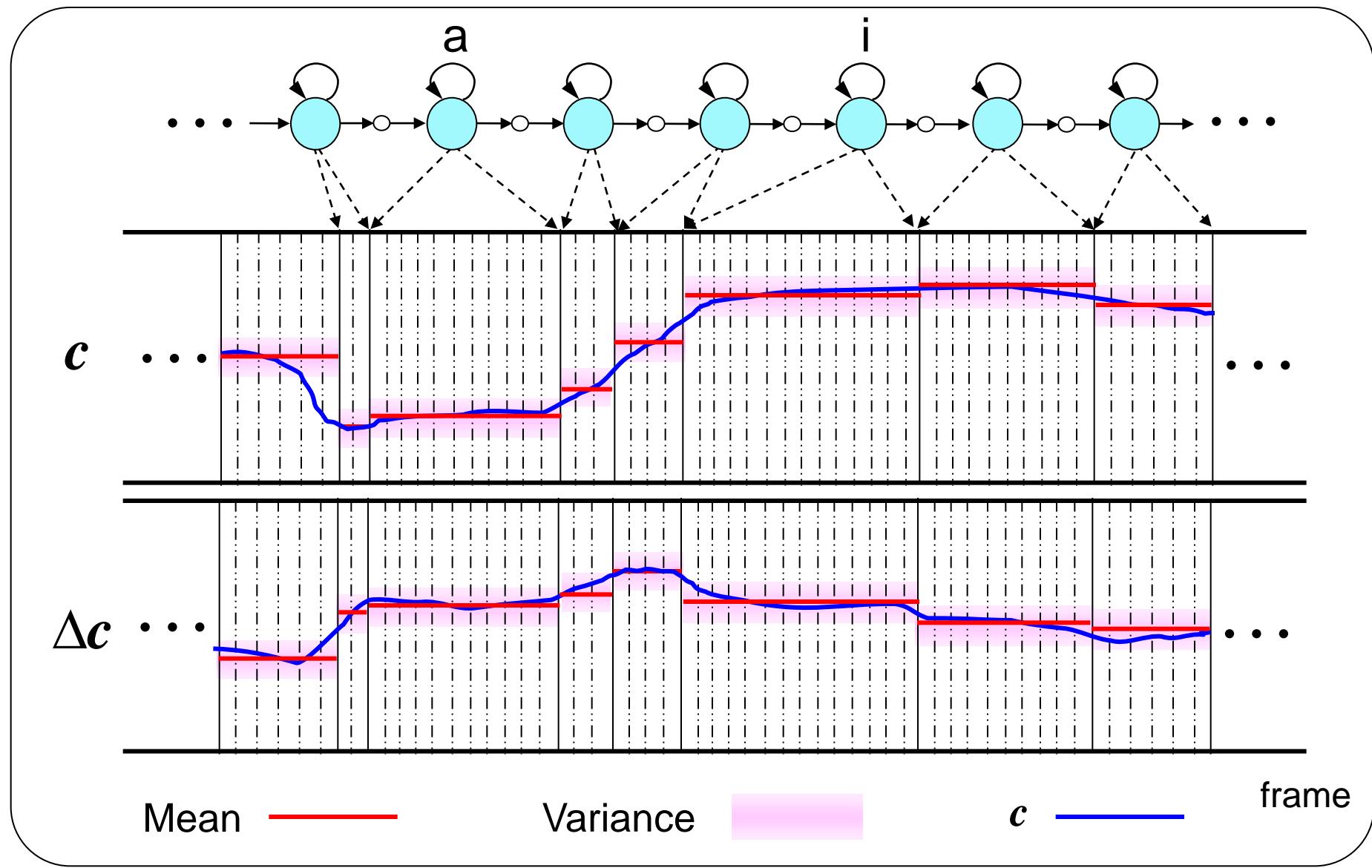
$$\mathbf{W}^\top \boldsymbol{\Sigma}_{\hat{\mathbf{q}}}^{-1} \mathbf{W} \mathbf{c} = \mathbf{W}^\top \boldsymbol{\Sigma}_{\hat{\mathbf{q}}}^{-1} \boldsymbol{\mu}_{\hat{\mathbf{q}}}$$

where

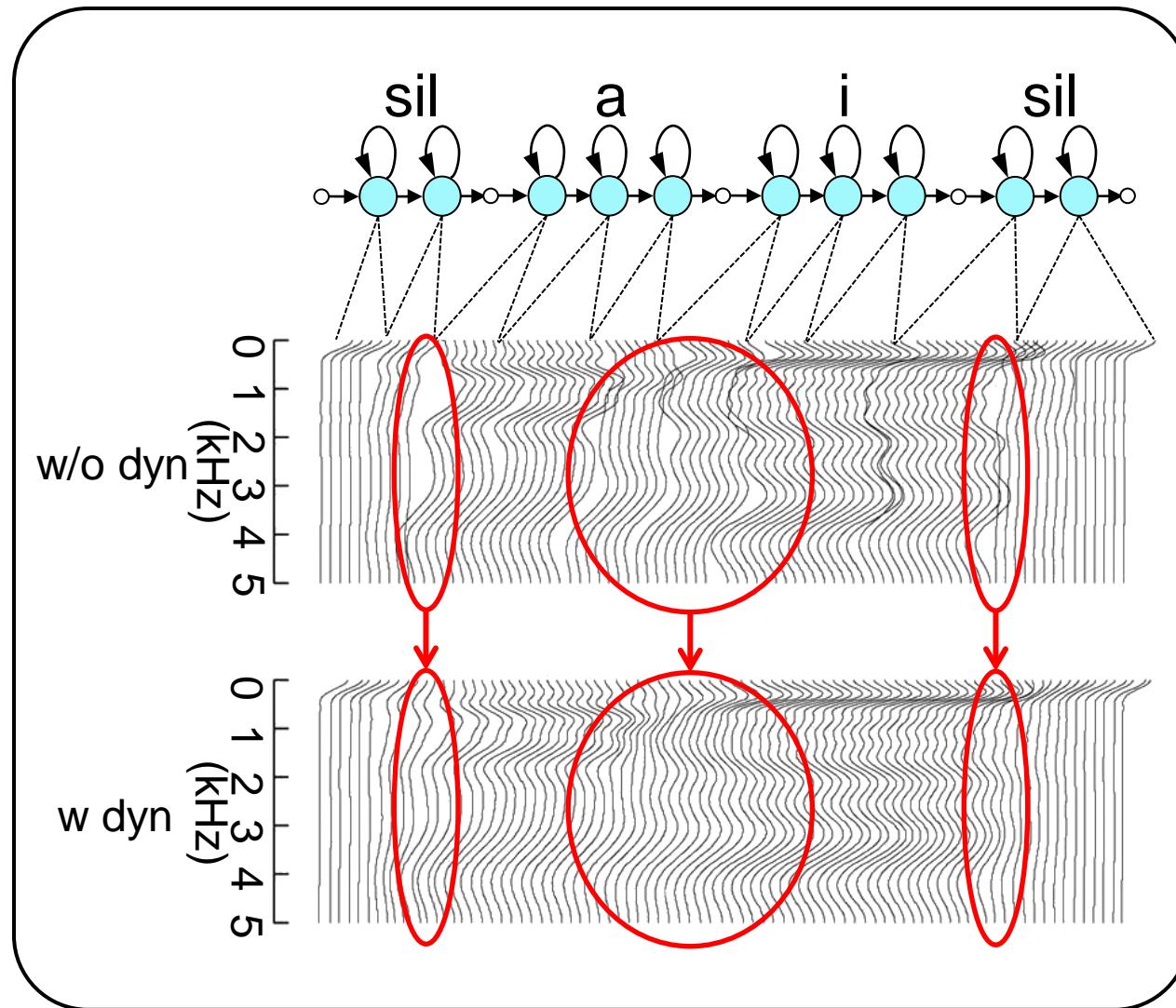
$$\mathbf{c} = [\mathbf{c}_1^\top, \mathbf{c}_2^\top, \dots, \mathbf{c}_T^\top]^\top$$

$$\boldsymbol{\mu}_{\hat{\mathbf{q}}} = [\boldsymbol{\mu}_{\hat{\mathbf{q}}_1}^\top, \boldsymbol{\mu}_{\hat{\mathbf{q}}_2}^\top, \dots, \boldsymbol{\mu}_{\hat{\mathbf{q}}_T}^\top]^\top$$

Generated speech parameter trajectory



Generated spectra



Spectra changing smoothly between phonemes

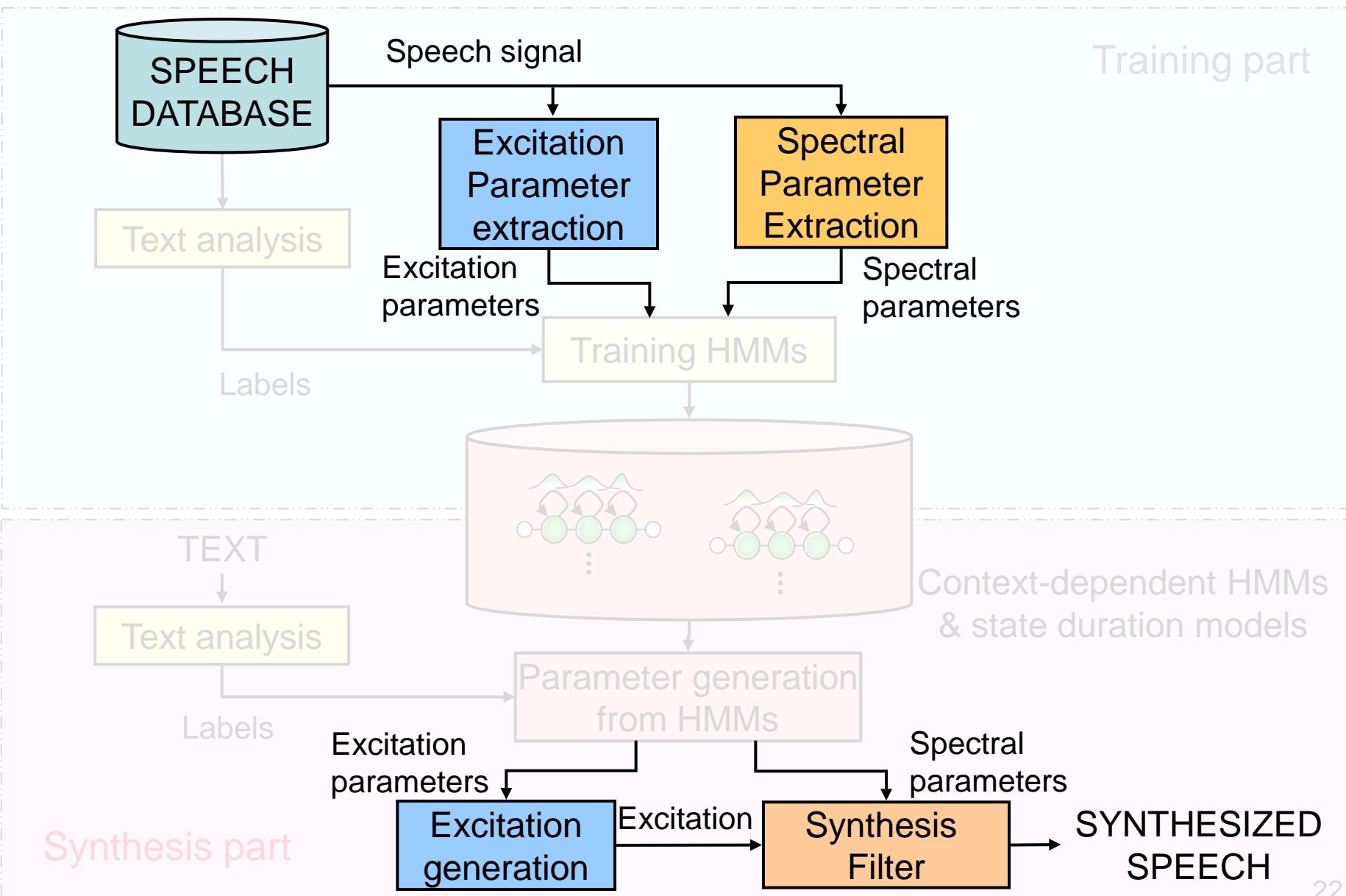
Trajectory HMM

$P(\mathbf{o} | \mathbf{l}, \lambda) = P(W\mathbf{c} | \mathbf{l}, \lambda)$ is not a proper distribution of \mathbf{c}

	Conventional HMM	Trajectory HMM
Training	$P(O L, \lambda)$	$P(C L, \lambda)$
Synthesis	$P(o l, \hat{\lambda}) _{o=Wc}$	$P(c l, \hat{\lambda})$

Solve inconsistency between training & synthesis
⇒ improving the model accuracy

HMM-based speech synthesis system

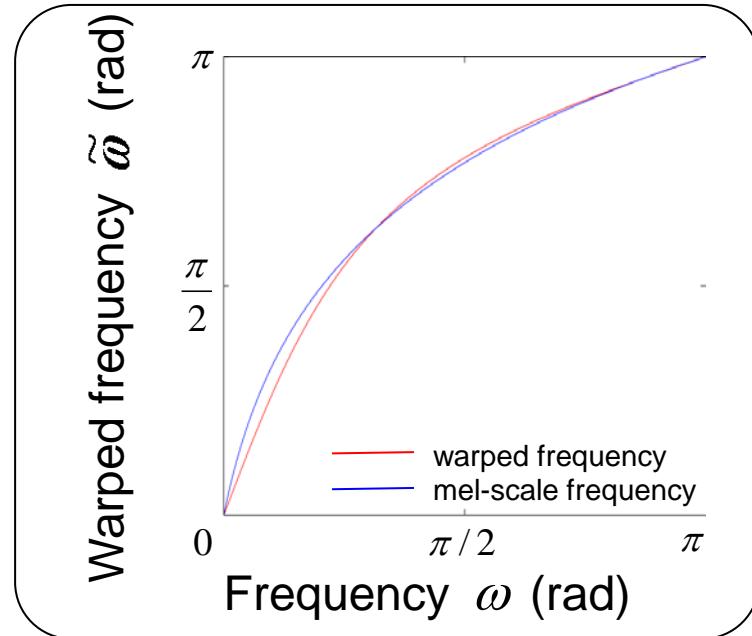


ML estimation of spectral parameter

Mel-cepstral representation of speech spectra

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

$$\tilde{z}^{-1} = \frac{z^{-1} - \alpha}{1 - \alpha z^{-1}} = e^{-j\tilde{\omega}}$$

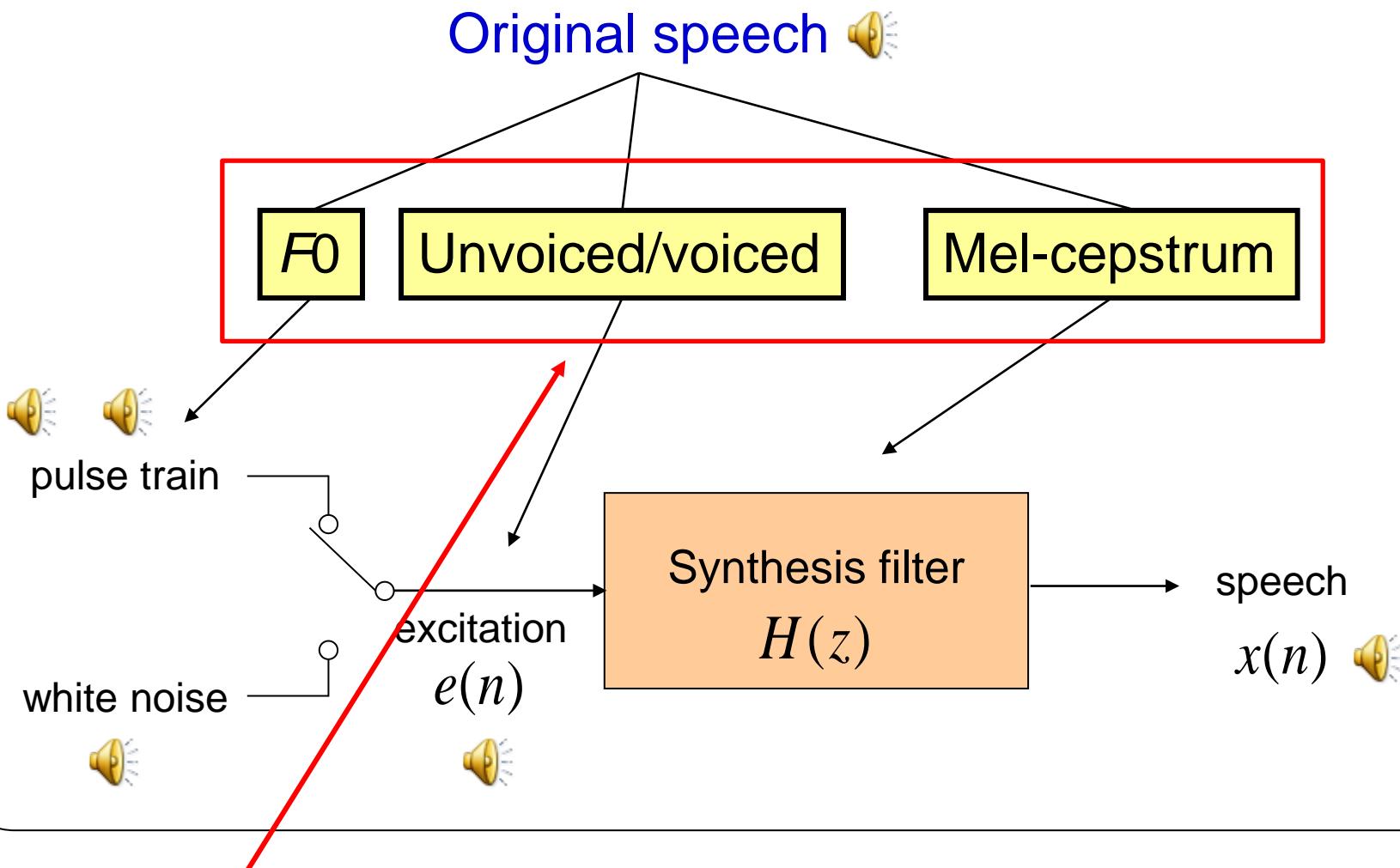


ML-estimation of mel-cepstrum

$$\mathbf{c} = \arg \max_c p(\mathbf{x} | \mathbf{c})$$

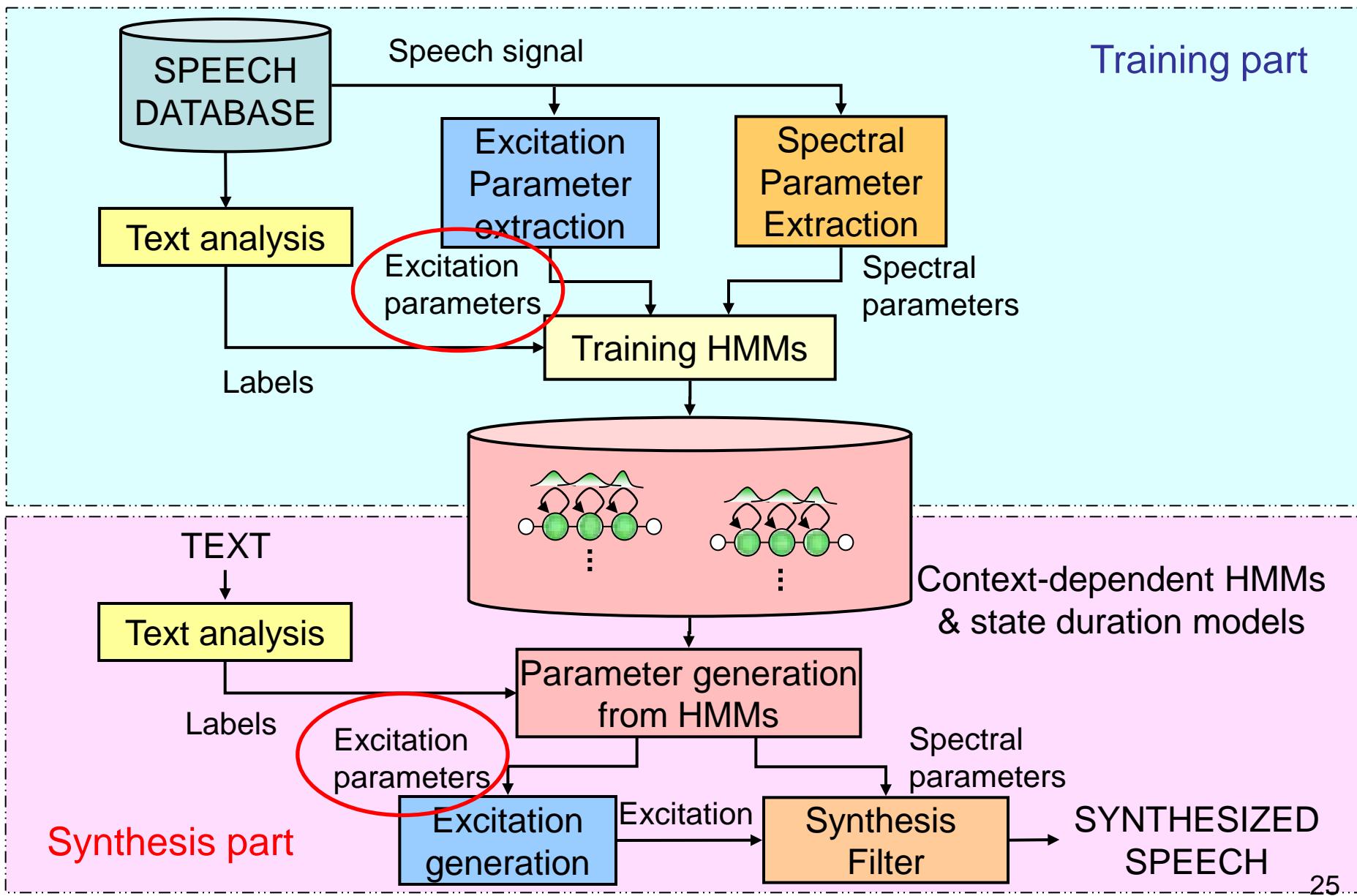
\mathbf{x} : speech waveform (Gaussian process)
 \mathbf{c} : mel-cepstrum

Overview of speech vocoding

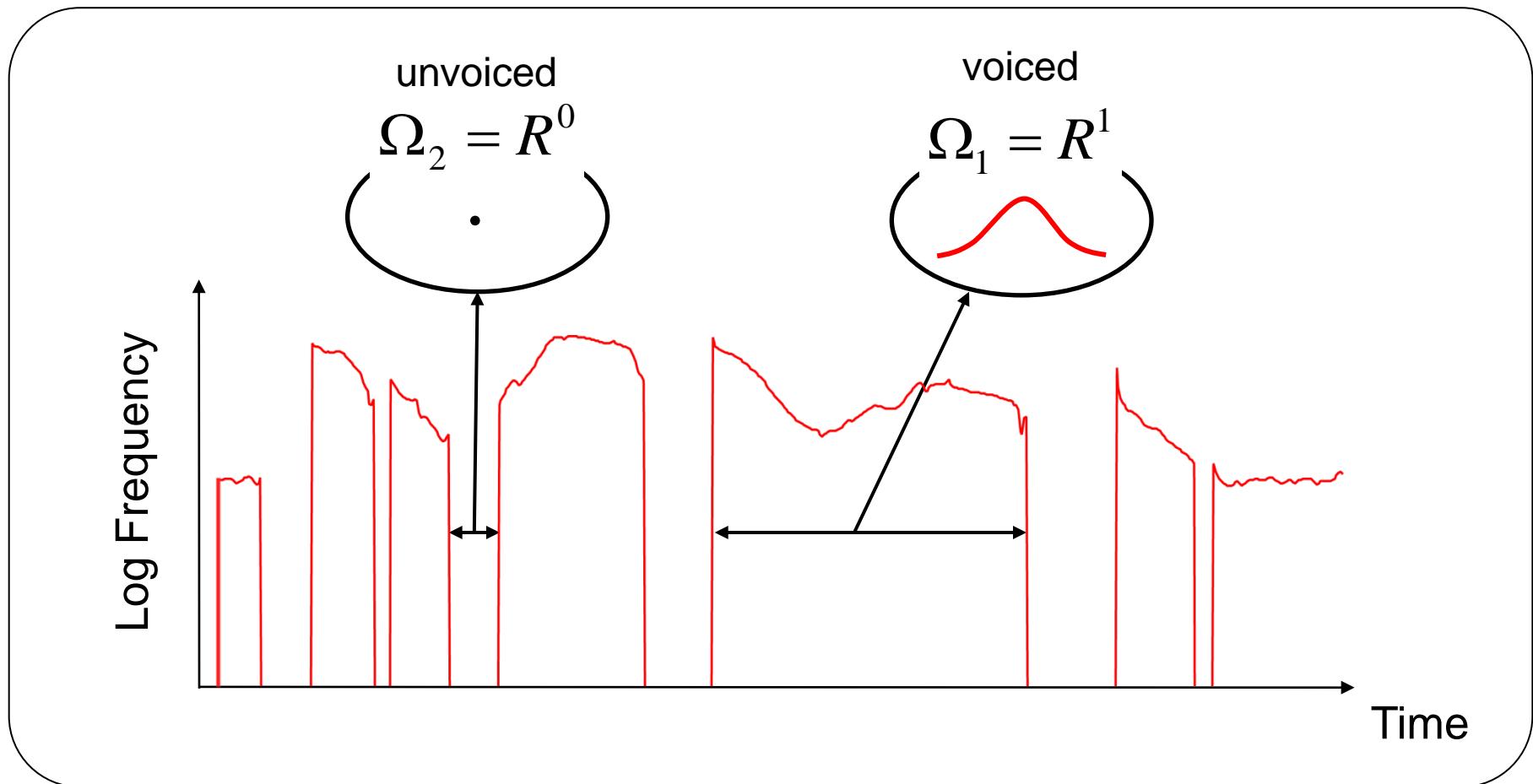


These speech parameter modeled by HMM

HMM-based speech synthesis system

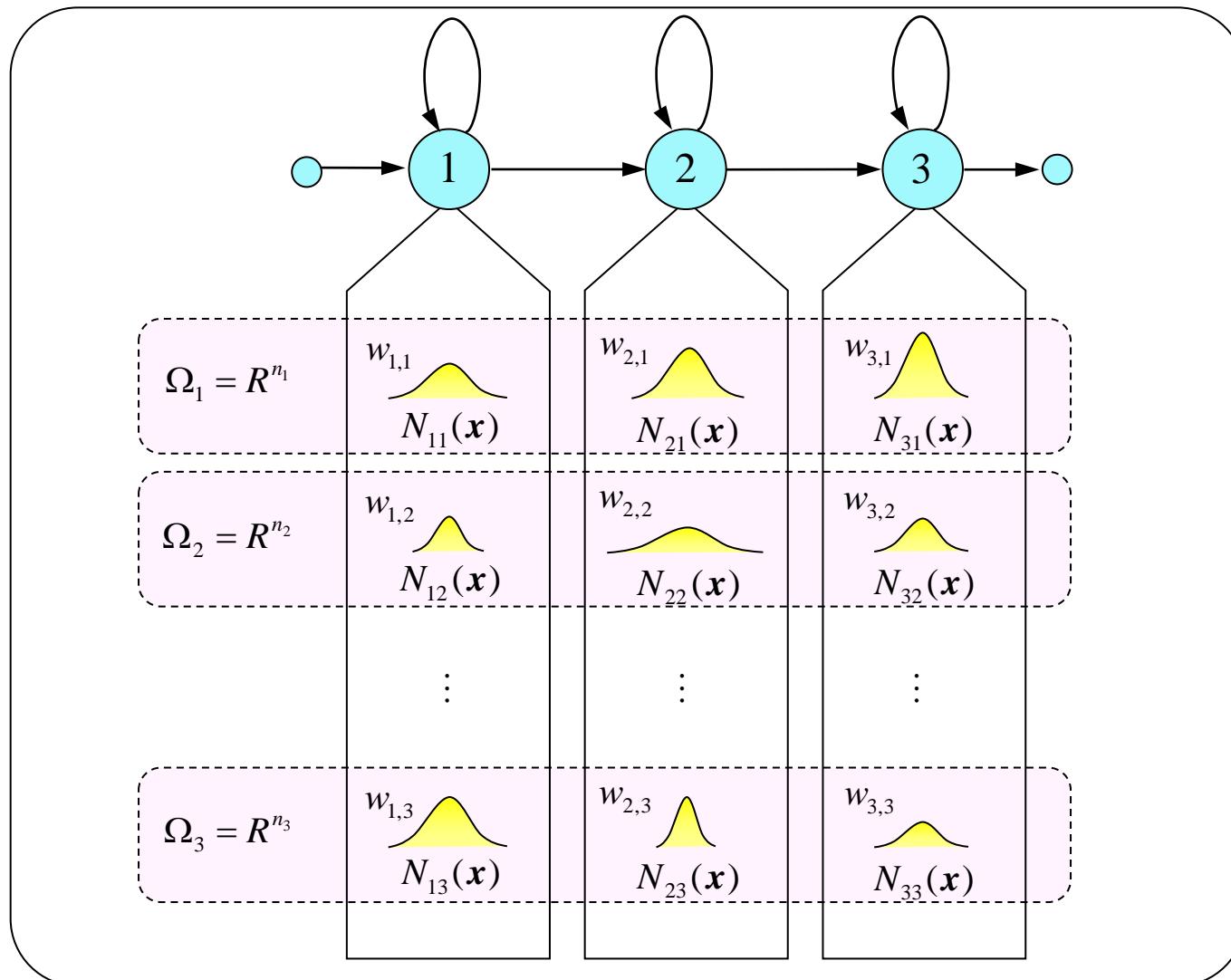


Observation of F0

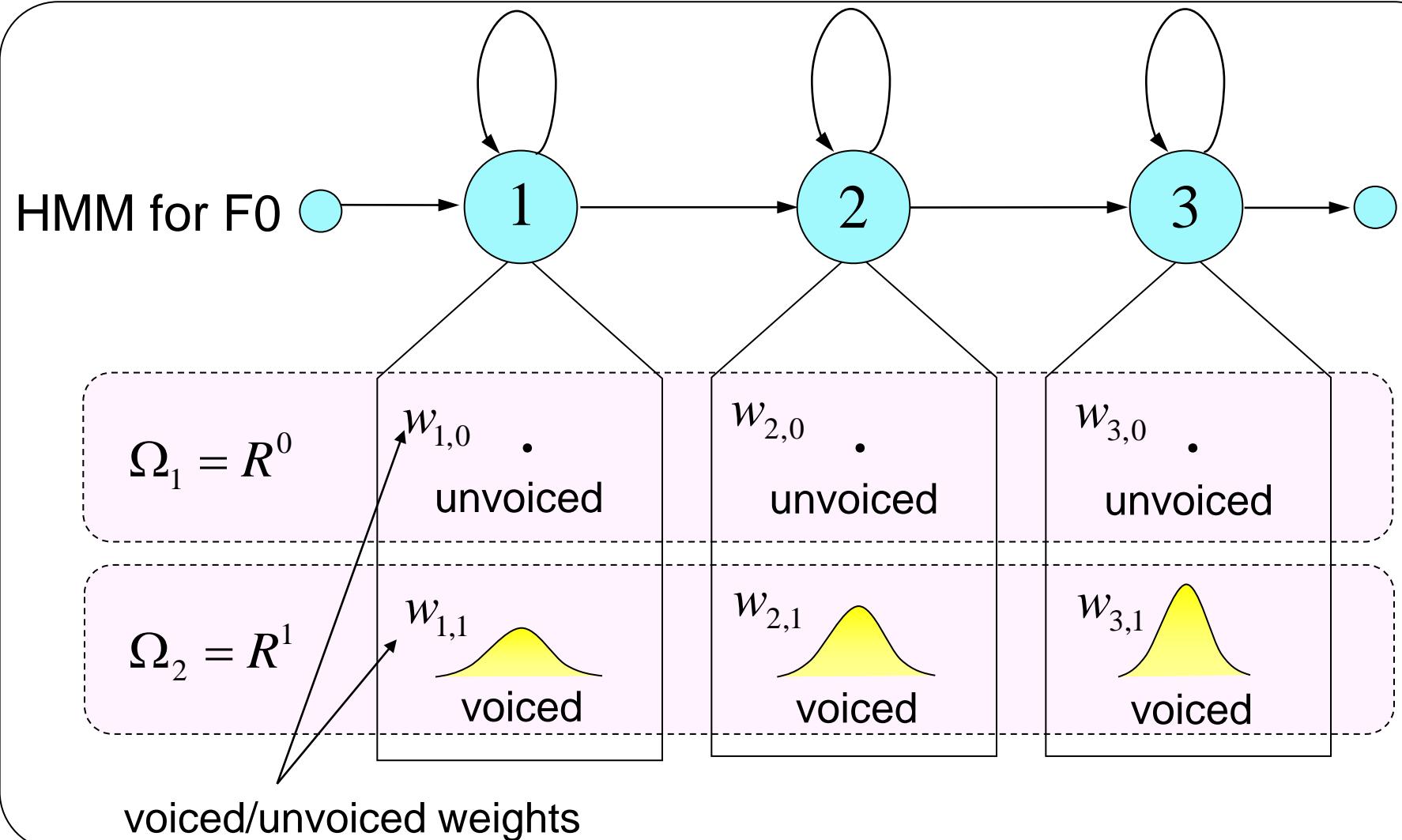


Unable to model by continuous or discrete distributions
⇒ Multi-space distribution HMM (MSD-HMM)

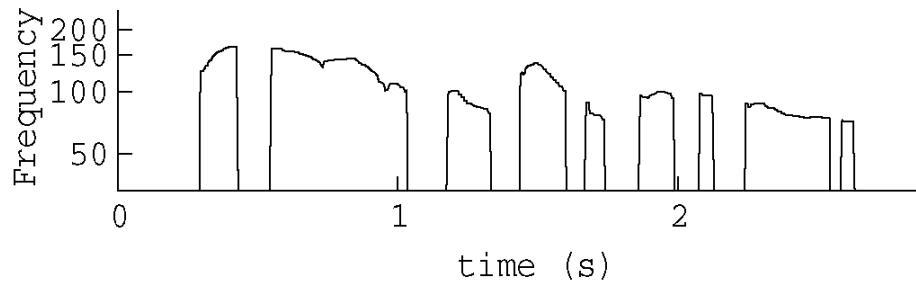
Structure of MSD-HMM



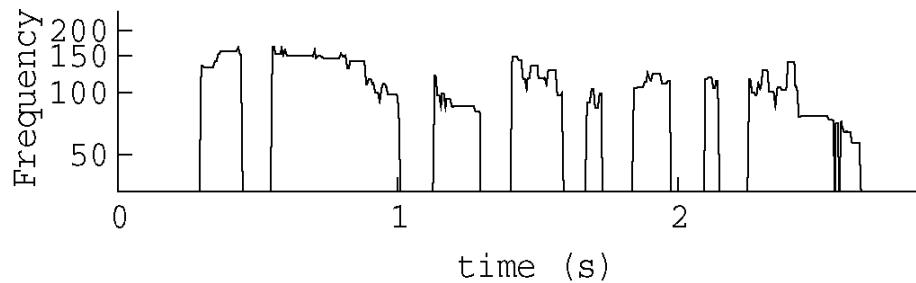
MSD-HMM for F0 modeling



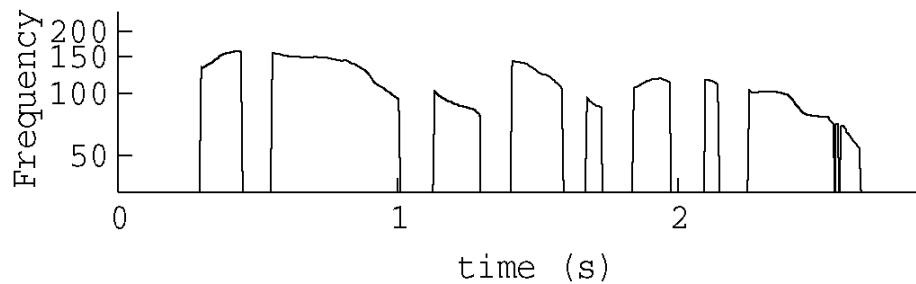
Generated F0



natural speech



without dynamic features



with dynamic features ($\Delta + \Delta^2$)

Speech samples

		Mel-cepstrum	
		w/ dyn.	w/o dyn.
log F0	w/ dyn.		
	w/o dyn.		

Inclusion of speech analysis & waveform reconstruction

Problem of statistical parametric speech synthesis

$$\hat{x} = \arg \max_x P(x | l, X, L)$$

$$= \arg \max_x \sum_c \sum_C \int \underbrace{P(x | c)}_{\text{Waveform reconstruction}} \underbrace{P(c | l, \lambda)}_{\text{Speech parameter generation}}$$

$$\times \underbrace{P(\lambda | C, L) P(C | X)}_{\text{HMM parameter estimation}} d\lambda \underbrace{\quad}_{\text{Speech analysis}}$$

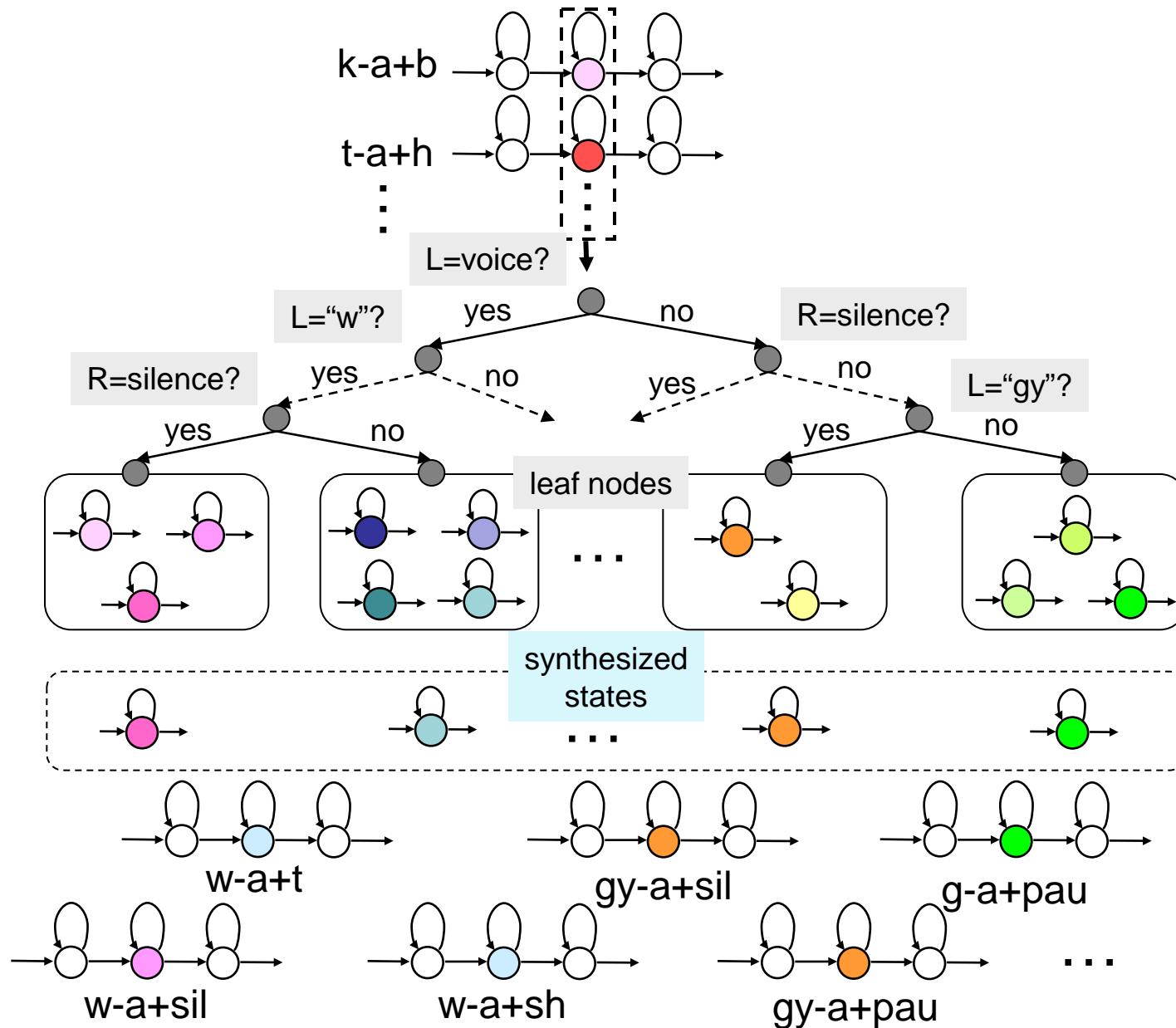
c and C consist of mel-cepstrum & F0 parameters

Context clustering factors

- {preceding, succeeding} two phonemes
 - Current phoneme
 - Position of current phoneme in current 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
 - Guess at 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
 - # of syllables in {preceding, current, succeeding} phrase
- ...

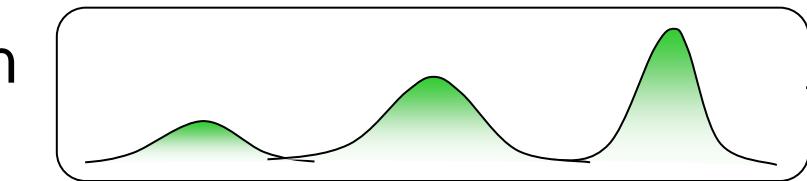
Vast # of combinations \Rightarrow Difficult to have possible models

Decision tree-based state clustering

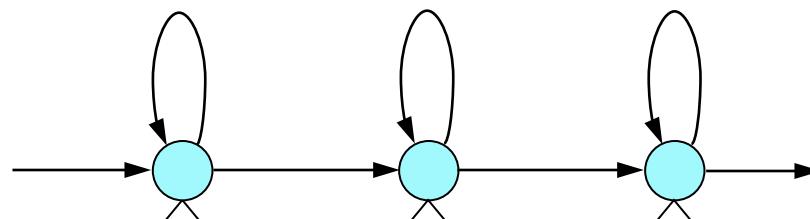


Stream-dependent tree-based clustering

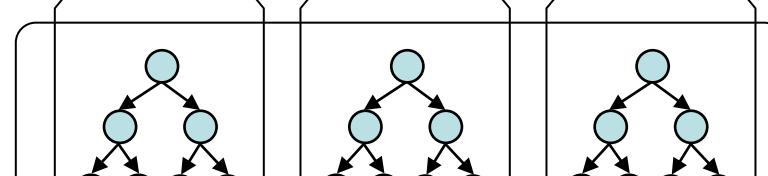
State duration model



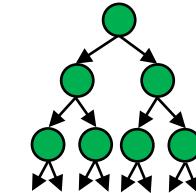
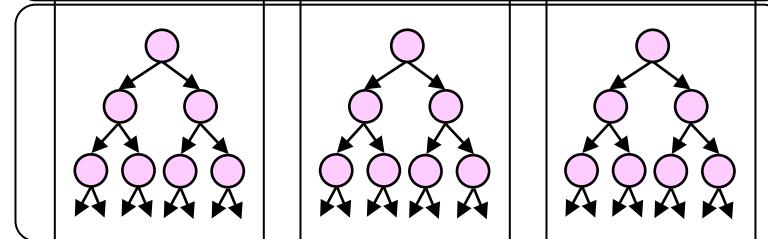
HMM



Decision trees
for
mel-cepstrum



Decision trees
for F0



Decision tree
for state dur.
models

Inclusion of model structure parameter

Problem of statistical parametric speech synthesis

$$\hat{\boldsymbol{o}} = \arg \max_{\boldsymbol{o}} P(\boldsymbol{o} | \boldsymbol{l}, \boldsymbol{O}, \boldsymbol{L})$$

$$= \arg \max_{\boldsymbol{o}} \sum_m \int \underbrace{P(\boldsymbol{o} | \boldsymbol{l}, \lambda, m)}_{\text{Generate coefficient}}$$

Generate coefficient

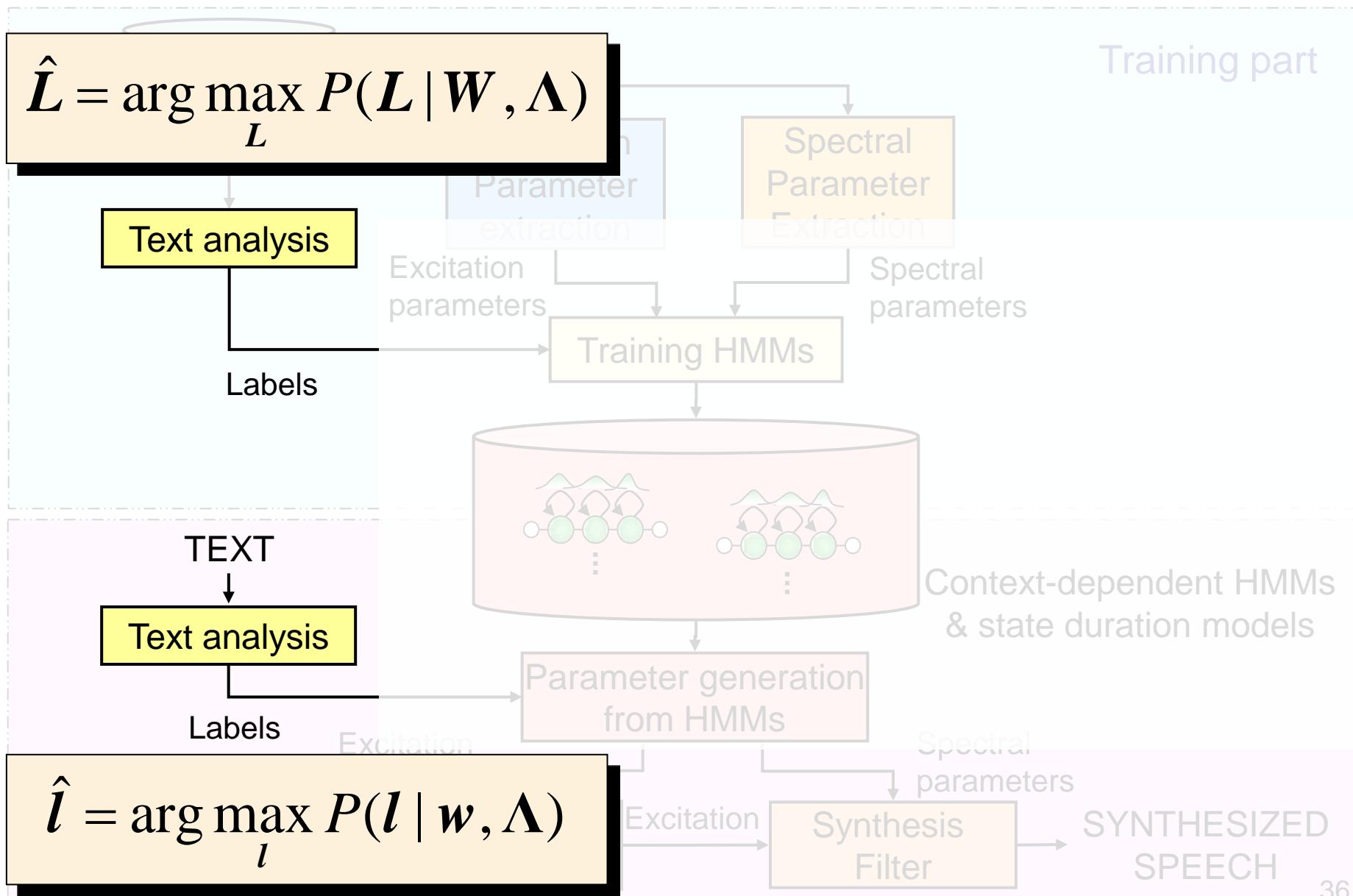
$$\times \underbrace{P(\lambda | m, \boldsymbol{O}, \boldsymbol{L})}_{\text{Posterior of model parameters}} \underbrace{P(m | \boldsymbol{O}, \boldsymbol{L}) d\lambda}_{\text{Posterior of model structure}}$$

Posterior of
model parameters

Posterior of
model structure

Usually fixed m is used

HMM-based speech synthesis system



Inclusion of text analysis

Problem of statistical parametric speech synthesis

$$\hat{\boldsymbol{o}} = \arg \max_{\boldsymbol{o}} P(\boldsymbol{o} | \boldsymbol{w}, \boldsymbol{O}, \boldsymbol{W})$$

$$= \arg \max_{\boldsymbol{o}} \sum_l \sum_L \int \int P(\boldsymbol{o} | \boldsymbol{l}, \boldsymbol{\lambda}) P(\boldsymbol{l} | \boldsymbol{w}, \boldsymbol{\Lambda})$$

Speech parameter generation

Text processing

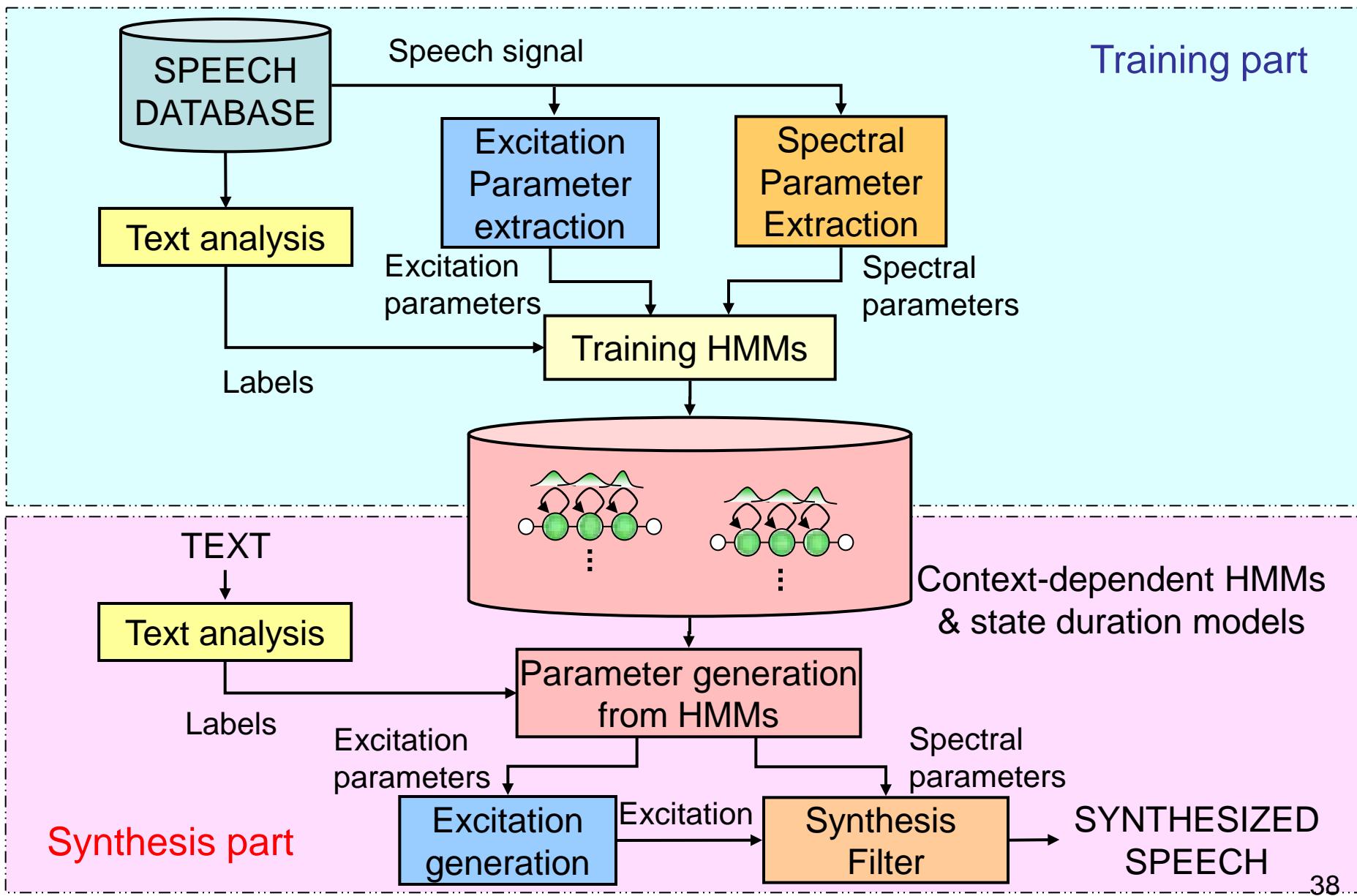
$$\times P(\boldsymbol{\lambda} | \boldsymbol{O}, \boldsymbol{L}) P(\boldsymbol{L} | \boldsymbol{W}, \boldsymbol{\Lambda}) P(\boldsymbol{\Lambda}) d\boldsymbol{\lambda} d\boldsymbol{\Lambda}$$

HMM parameter estimation

Text processing

Prior

HMM-based speech synthesis system



Inclusion of all components

Problem of statistical parametric speech synthesis

$$\hat{x} = \arg \max_x P(x | w, X, W)$$
$$= \arg \max_x \sum_{c,C} \sum_{l,L} \sum_m \int \int \underbrace{P(x | c)}_{\text{Waveform generation}} \underbrace{P(c | l, \lambda, m)}_{\text{Parameter generation}} \underbrace{P(l | w, \Lambda)}_{\text{Text processing}}$$
$$\times \underbrace{P(\lambda | m, C, L)}_{\text{Posterior of model parameter}} \underbrace{P(m | C, L)}_{\text{Posterior of model structure}}$$
$$\times \underbrace{P(L | W, \Lambda)}_{\text{Text processing}} \underbrace{P(C | X)}_{\text{Speech analysis}} \underbrace{P(\Lambda)}_{\text{Prior}} d\lambda d\Lambda$$

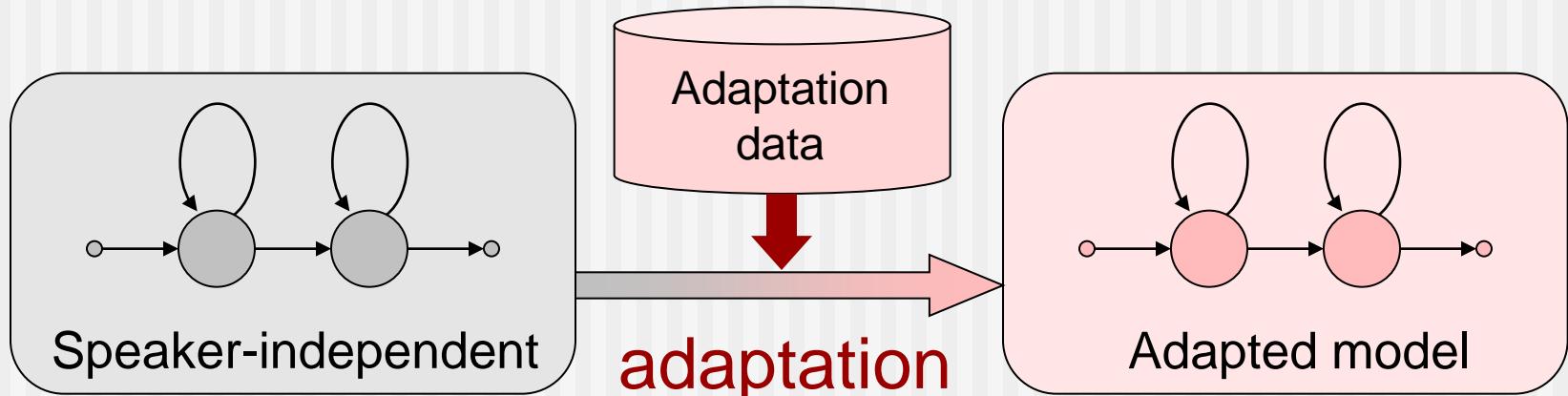
Emotional Speech Synthesis

text	neutral	angry
<p>「授業中に携帯いじってんじゃねえよ！ 電源切っとけ！」</p> <p>“Don’t touch your cell phone during a class! Turn off it!”</p>		
<p>「ミーティングには毎週参加しなさい！」</p> <p>“You must attend the weekly meeting!”</p>		

trained with 200 utterances

Speaker Adaptation (mimicking voices)

MLLR-based adaptation



- w/o adaptation (initial model)
- Adapted with 4 utterances
- Adapted with 50 utterances
- Speaker-dependent model

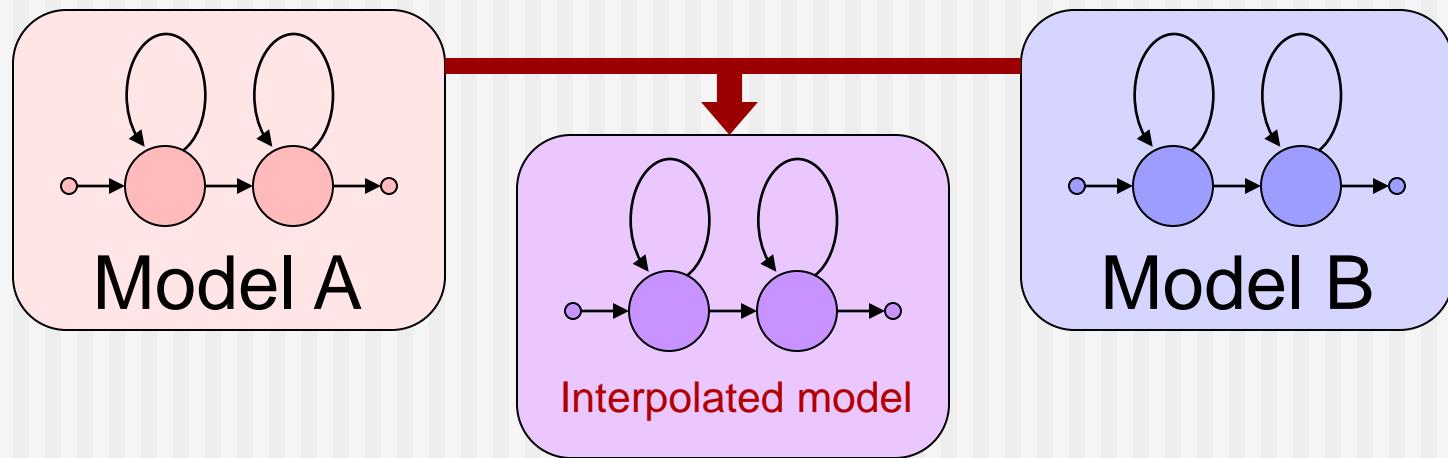


?



Speaker Interpolation (mixing voices)

Linear combination of two speaker-dependent models



A: 1.00

0.75

0.50

0.25

0.00



B: 0.00

0.25

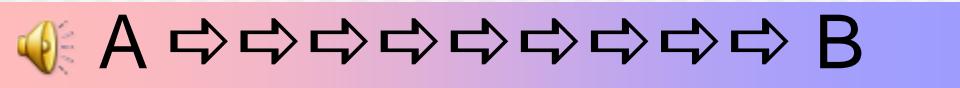
0.50

0.75

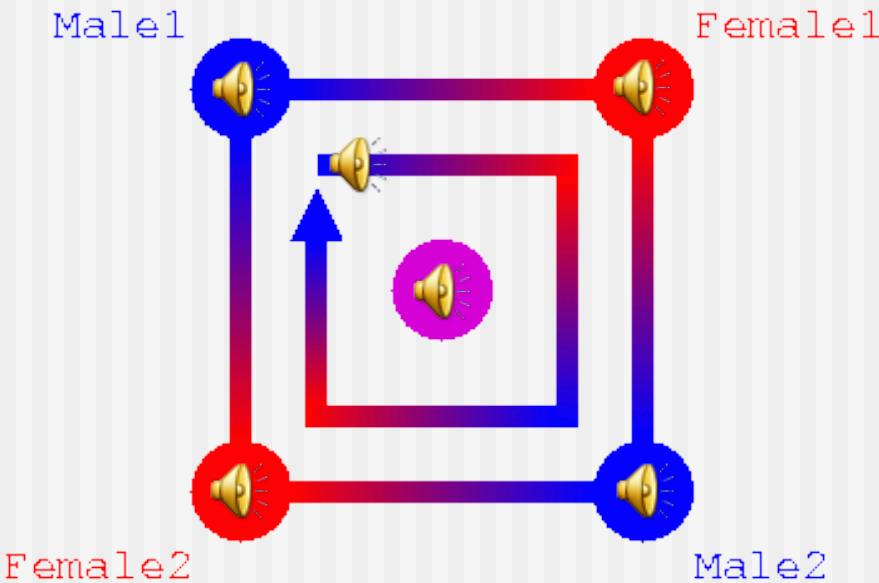
1.00

Voice Morphing

Two voices:

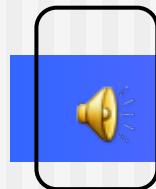


Four voices:



Interpolation of Speaking Styles

Base model A

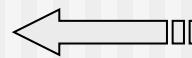


Interpolation

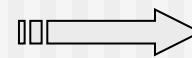
Base model B



extrapolation



Neutral

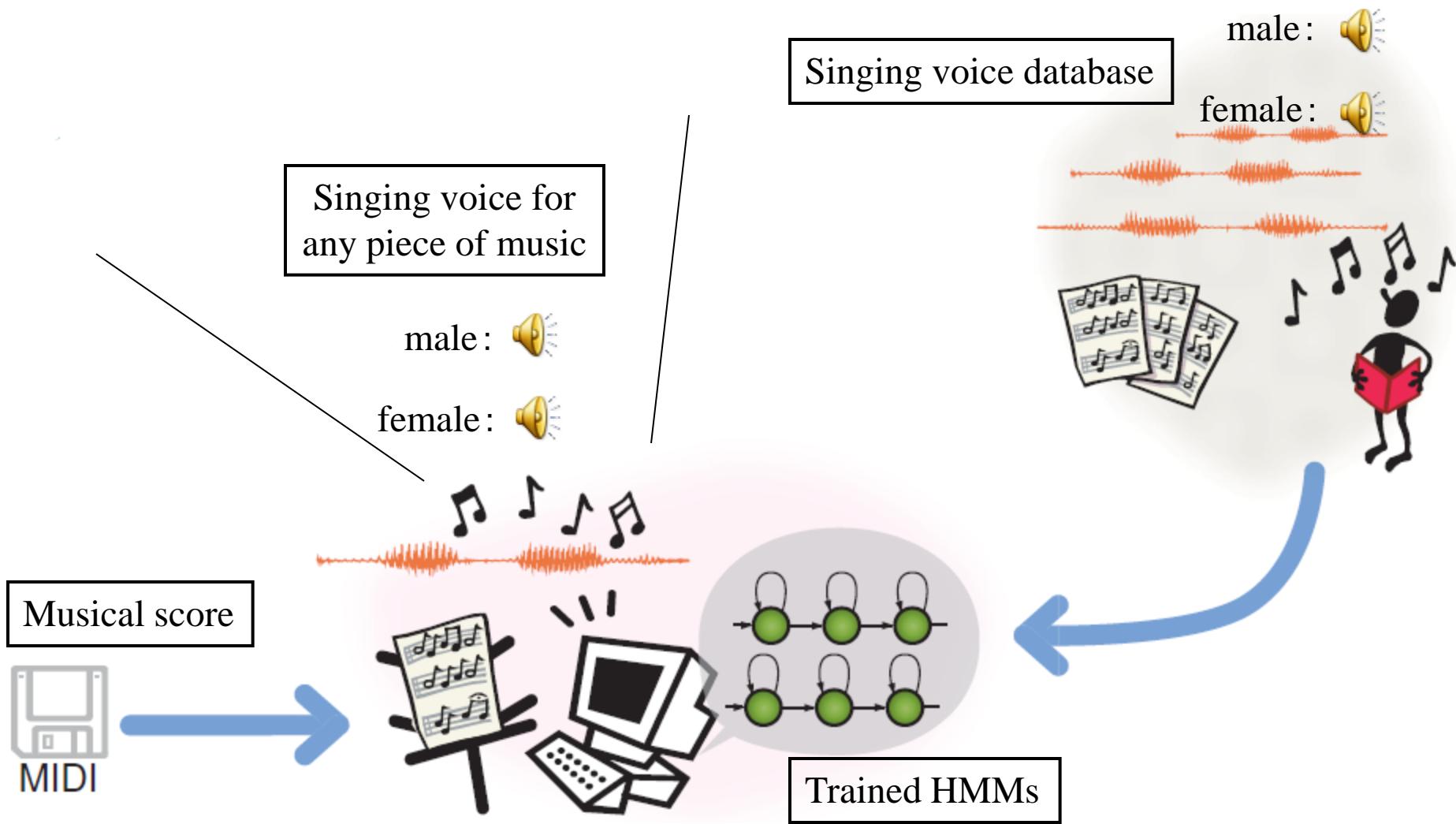


High Tension

Multilingual Speech Synthesis

- Japanese 🔊 🔊
- American English 🔊 🔊 🔊 🔊 🔊
- Chinese (Mandarin) (by ATR) 🔊
- Brazilian Portuguese (by Nitech, and UFRJ) 🔊
- European Portuguese 🔊 (by Nitech, Univ of Porto, and UFRJ)
- Slovenian 🔊 (by Bostjan Vesničer, University of Ljubljana, Slovenia)
- Swedish 🔊 🔊 (by Anders Lundgren, KTH, Sweden)
- German (by University of Bonn, and Nitech) 🔊
- Korean (by Sang-Jin Kim, ETRI, Korea) 🔊 🔊
- Hungarian (by Budapest University of Technology and Economics) 🔊
- Finish (by TKK, Finland) 🔊 🔊
- Baby English (by Univ of Edinburgh, UK) 🔊
- Polish, Slovak, Finnish, Arabic, Farsi, Polyglot, etc.

Singing Voice Synthesis



Conclusion

Statistical approach to speech synthesis

- Whole speech synthesis process is described as a statistical framework
- It gives a unified view and reveals what is correct and what is wrong
- Importance of the database

Future work

- Still we have many problems should be solved:
 - Speech waveform modeling
 - Combination with text processing part