

基于HMM的语音合成

HMM based Speech Synthesis

清华大学深圳研究生院

吴志勇

zywu@sz.tsinghua.edu.cn



Adapted from HTS Slides

HTS Slides
released by HTS Working Group
<http://hts.sp.nitech.ac.jp/>

Copyright (c) 1999 - 2011
Nagoya Institute of Technology
Department of Computer Science

Some rights reserved.

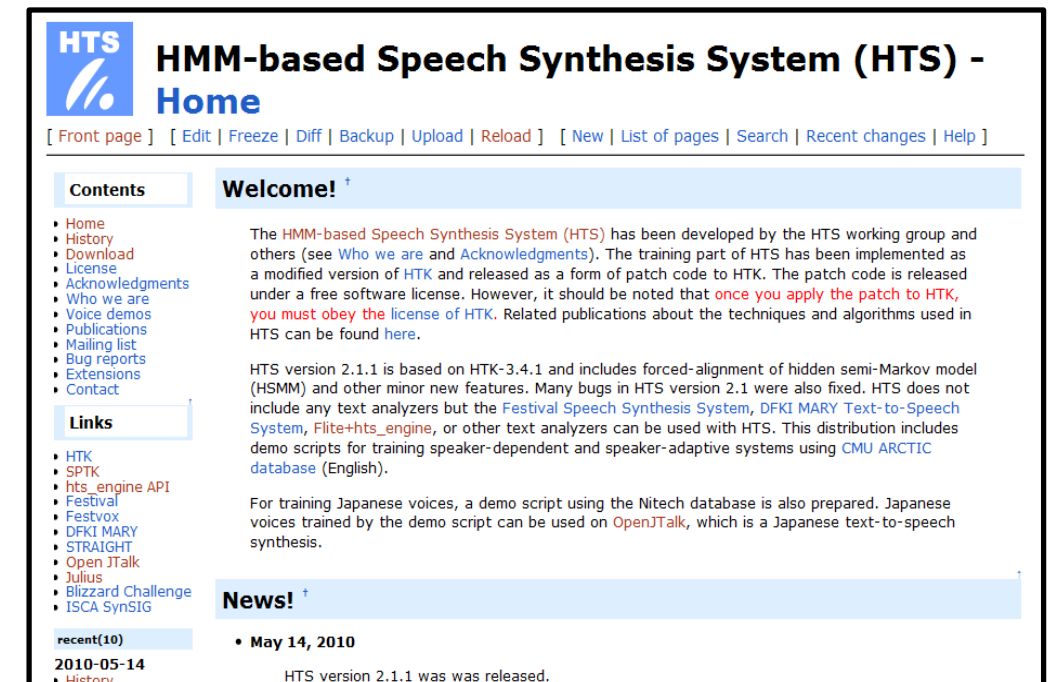
This work is licensed under the Creative Commons Attribution 3.0 license.
See <http://creativecommons.org/> for details.



HTS: HMM-base Speech Synthesis System

■ HMM-based Speech Synthesis System (HTS)

- ❑ Released as a form of patch code to HTK
 - Under the New and Simplified BSD license
 - ❑ Once you apply the patch to HTK, you must obey the license of HTK
- ❑ HTS-users mailing list
 - Over 500 posts per year
 - All posts are archived & searchable
 - Bug reports, Q&A, announce
- ❑ Becoming a research platform
 - Using by various organizations (e. g., Microsoft, IBM)

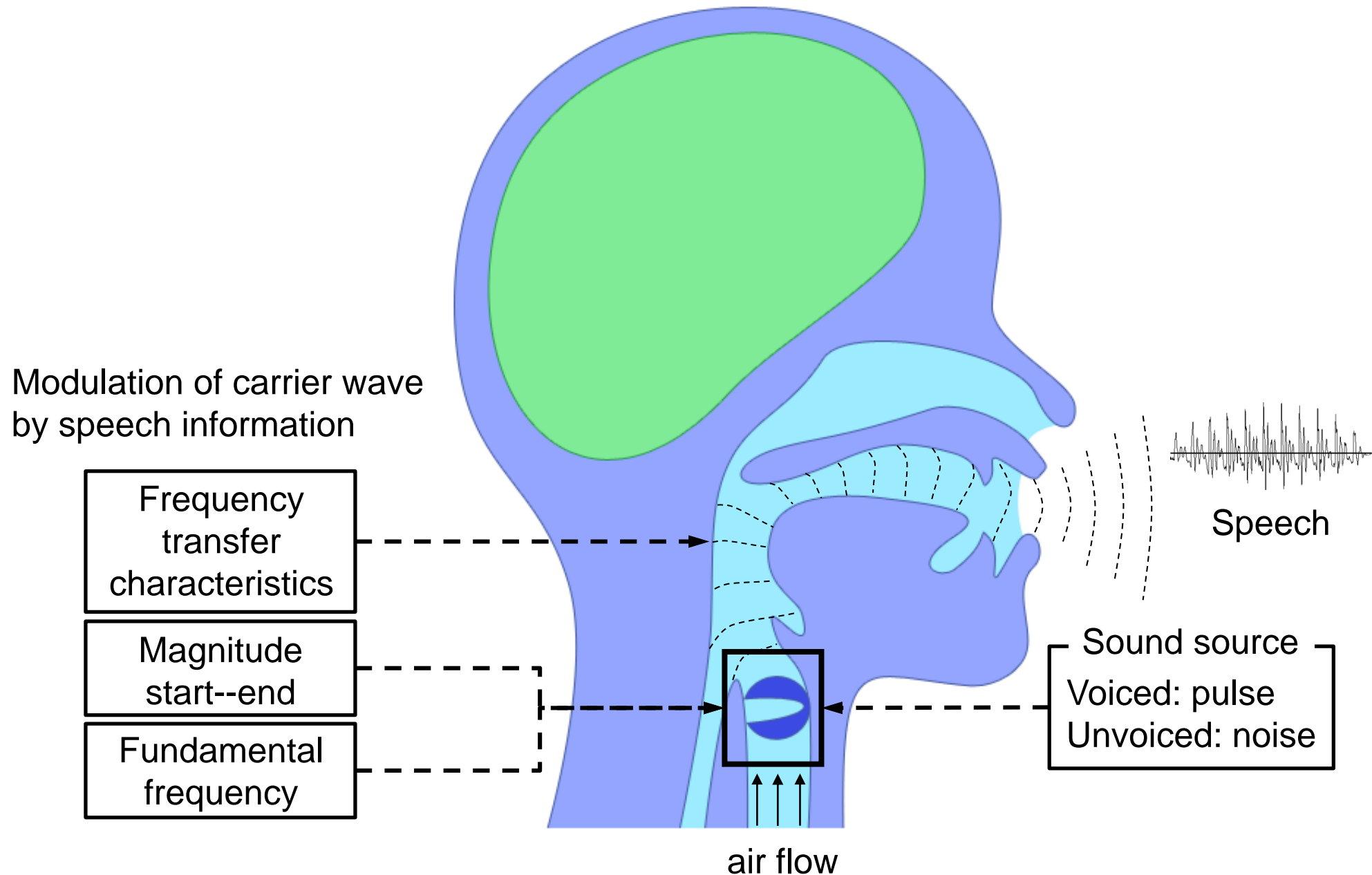


<http://hts.sp.nitech.ac.jp/>

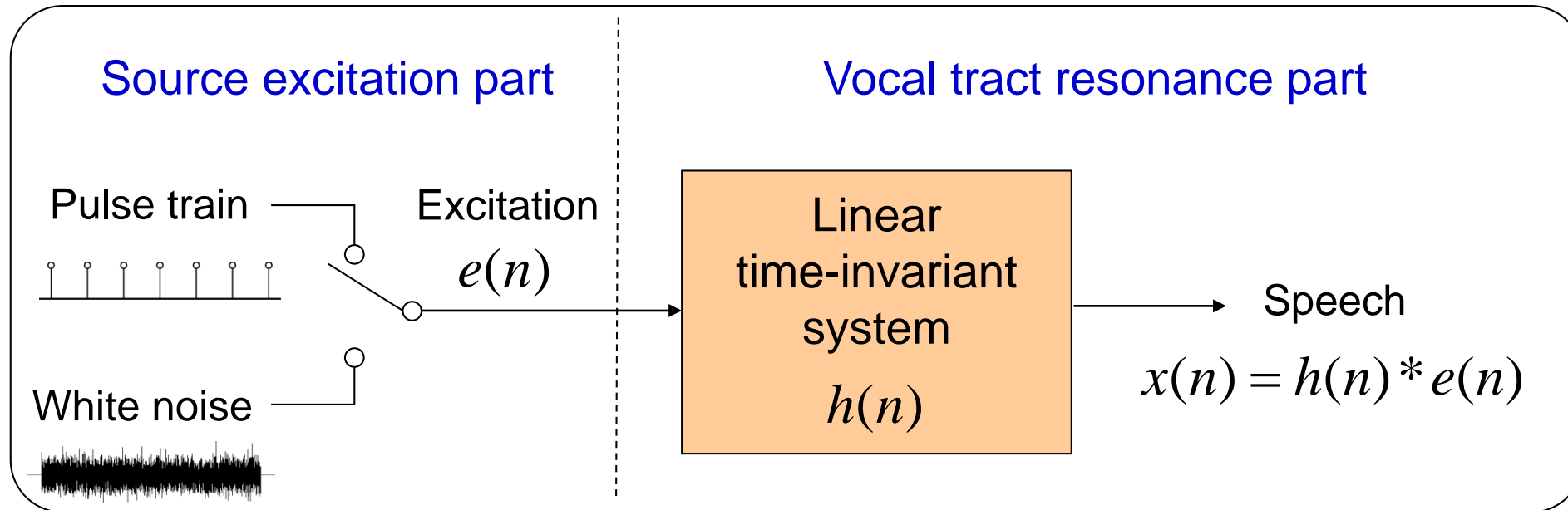
- **Speech vocoding: Source-filter model**
- **Speech parameter modeling and generation with HMM**
 - Overview of HMM framework
 - State duration modeling
 - Spectrum modeling
 - F0 modeling
 - Context clustering
- **Voice character controlling**
 - Adaptation (mimicking voices)
 - Interpolation (mixing voices)
- **Application**
 - Singing voice
 - Emotional voice
 - Audio-visual speech synthesis

- **Speech vocoding: Source-filter model**
- **Speech parameter modeling and generation with HMM**
 - Overview of HMM framework
 - State duration modeling
 - Spectrum modeling
 - F0 modeling
 - Context clustering
- **Voice character controlling**
 - Adaptation (mimicking voices)
 - Interpolation (mixing voices)
- **Application**
 - Singing voice
 - Emotional voice
 - Audio-visual speech synthesis

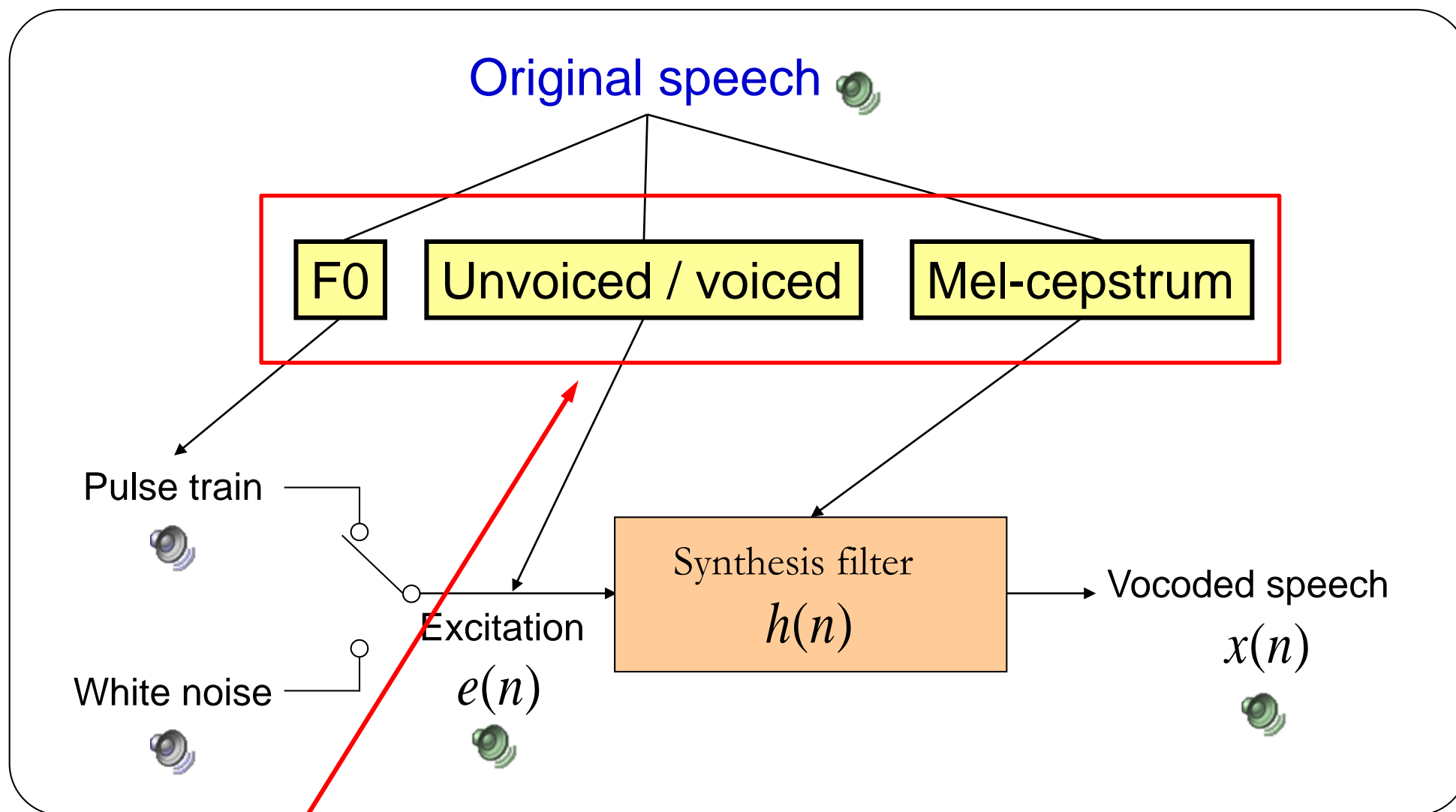
Speech Production Mechanism



Source-Filter Model

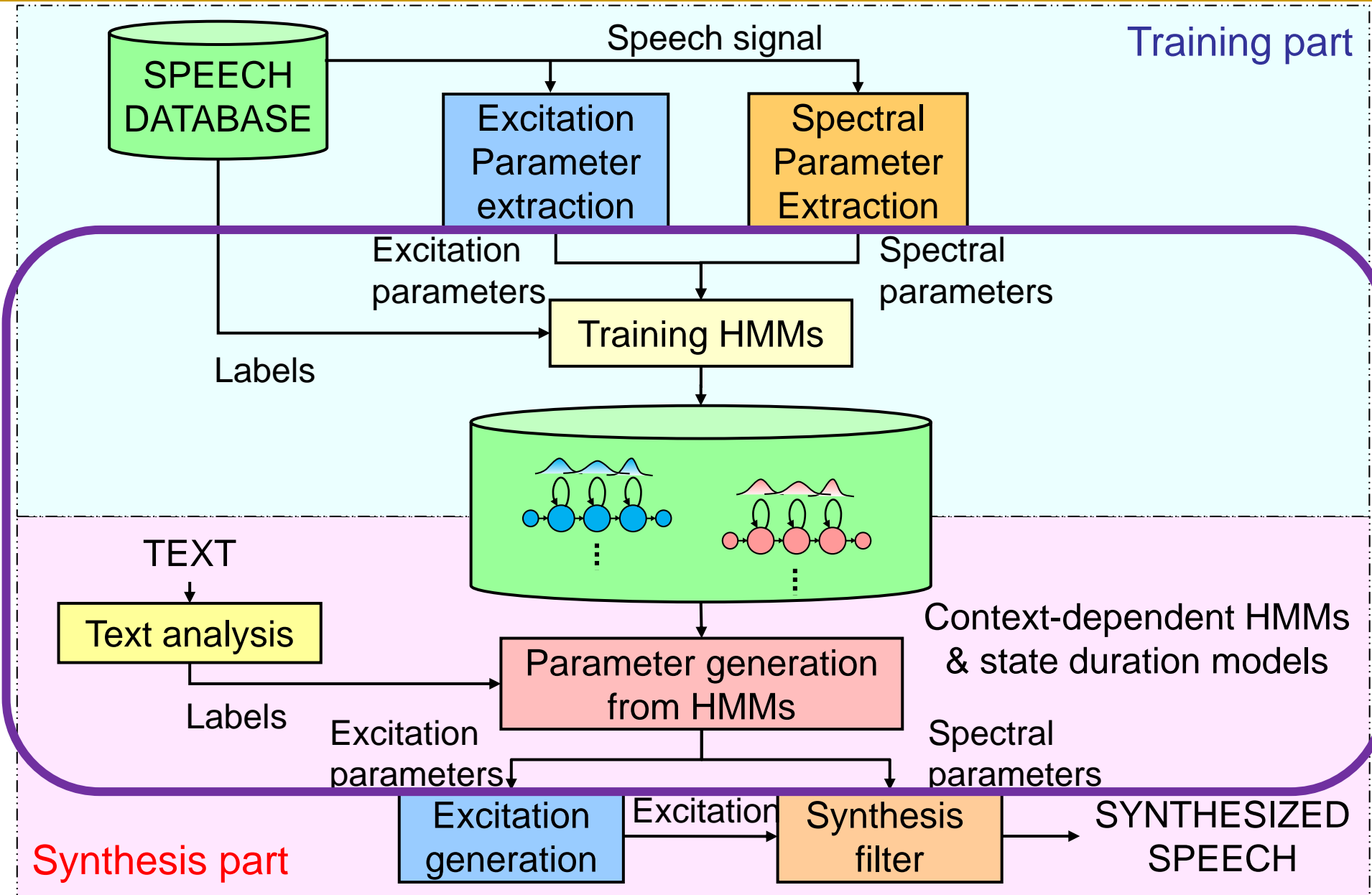



Overview of Speech Vocoding



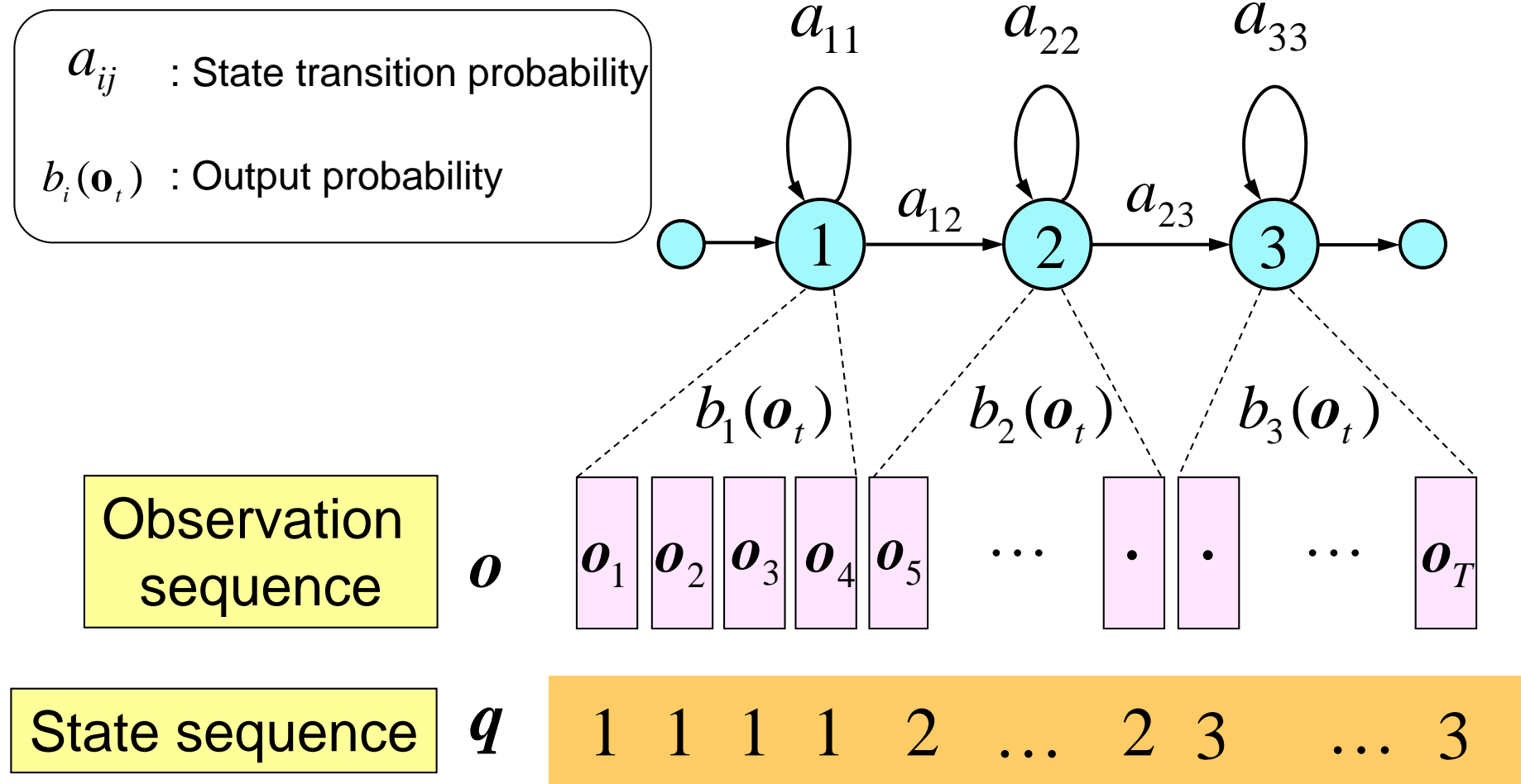
These speech parameters are to be modeled by HMM

HMM-based Speech Synthesis System



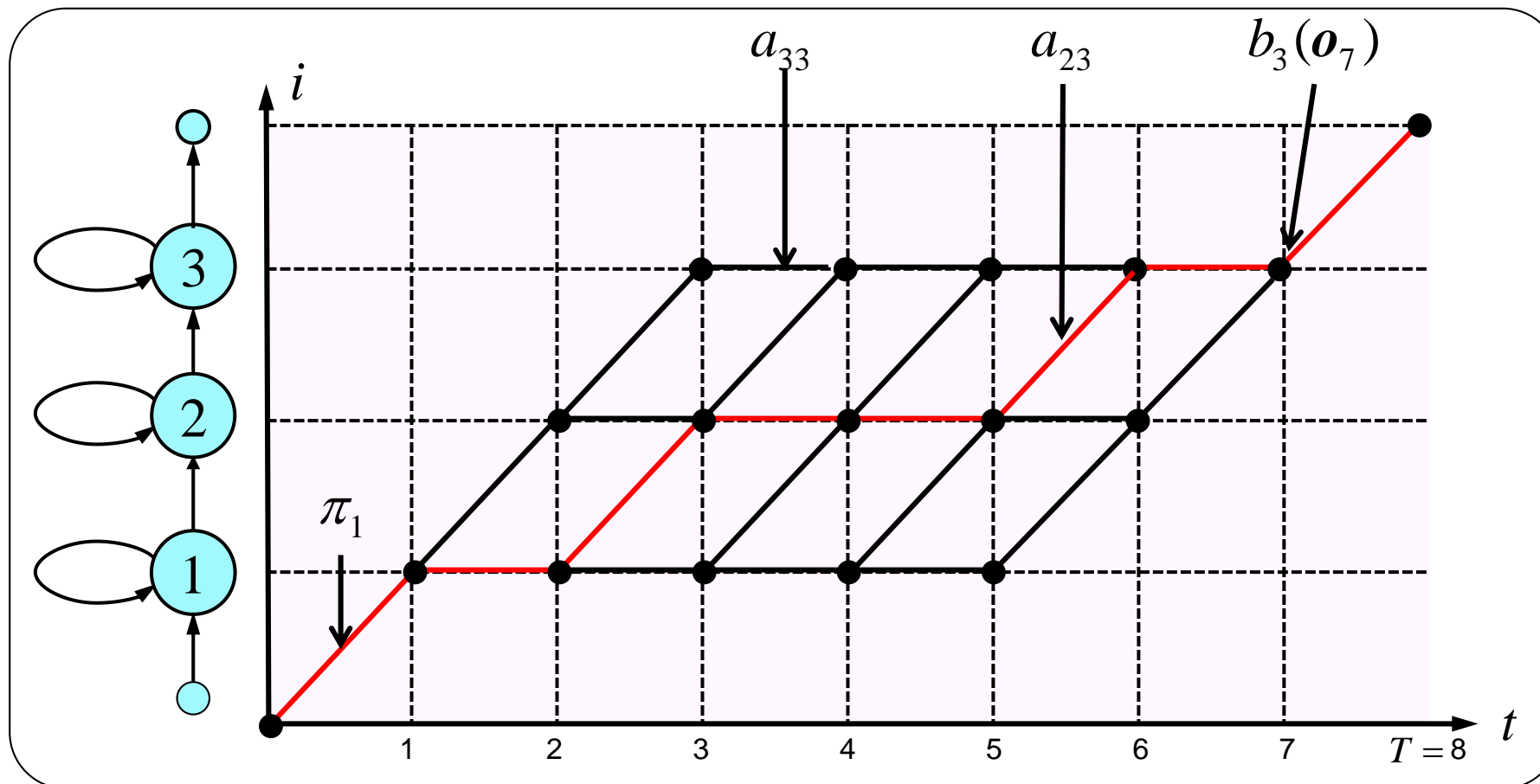
- Speech vocoding: Source-filter model
- **Speech parameter modeling and generation with HMM**
 - Overview of HMM framework 
 - State duration modeling
 - Spectrum modeling
 - F0 modeling
 - Context clustering
- Voice character controlling
 - Adaptation (mimicking voices)
 - Interpolation (mixing voices)
- Application
 - Singing voice
 - Emotional voice
 - Audio-visual speech synthesis

Hidden Markov Model (HMM)



The Markov chain whose state sequence is unknown
 \Rightarrow Estimating state sequence by the observation

Output Probability of HMM



Likelihood function

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

Speech Parameter Generation Algorithm [Tokuda; '00]

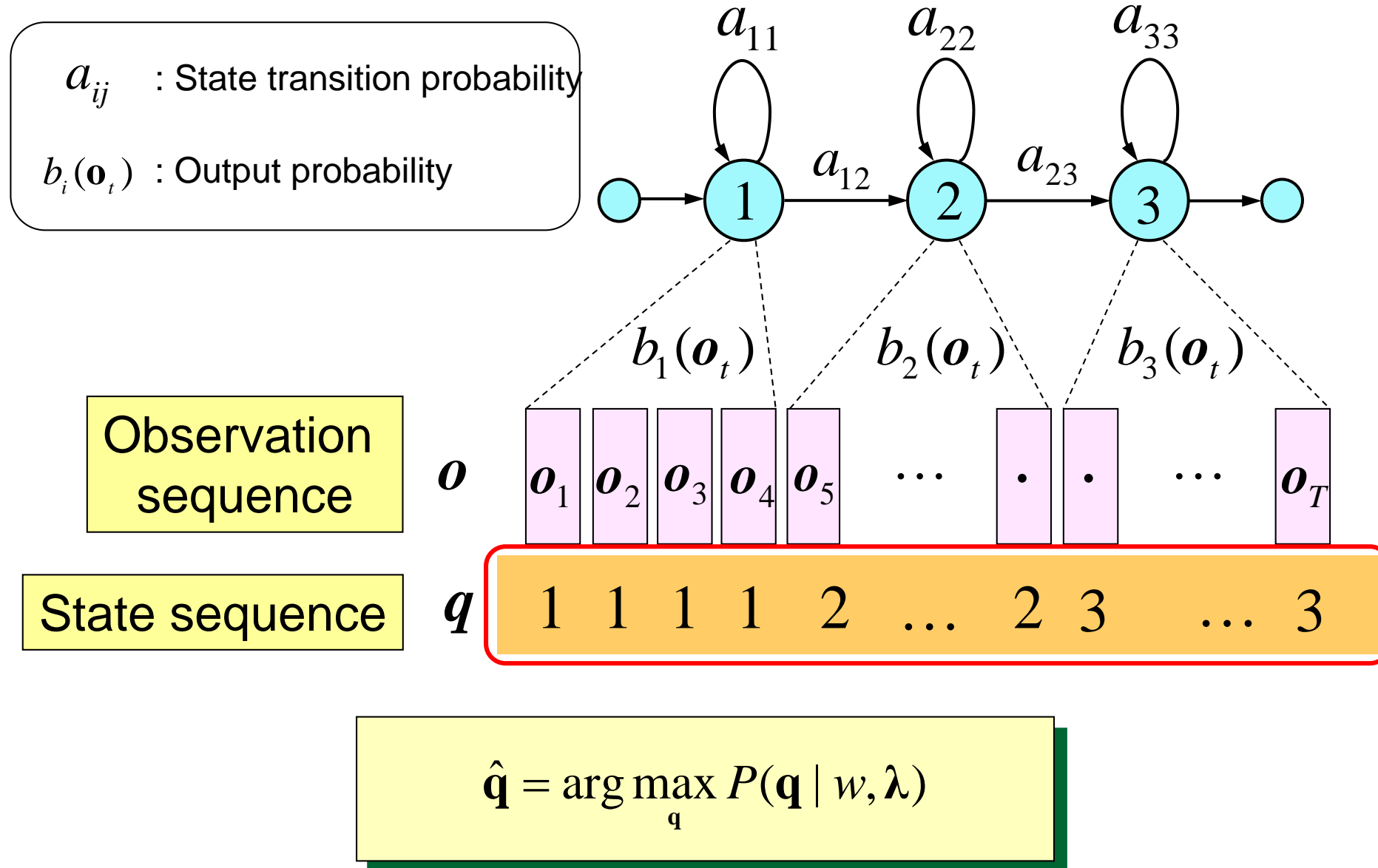
For a given HMM λ , determine a speech parameter vector sequence $\mathbf{o} = [\mathbf{o}_1^\top, \mathbf{o}_2^\top, \dots, \mathbf{o}_T^\top]^\top$ which maximizes


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



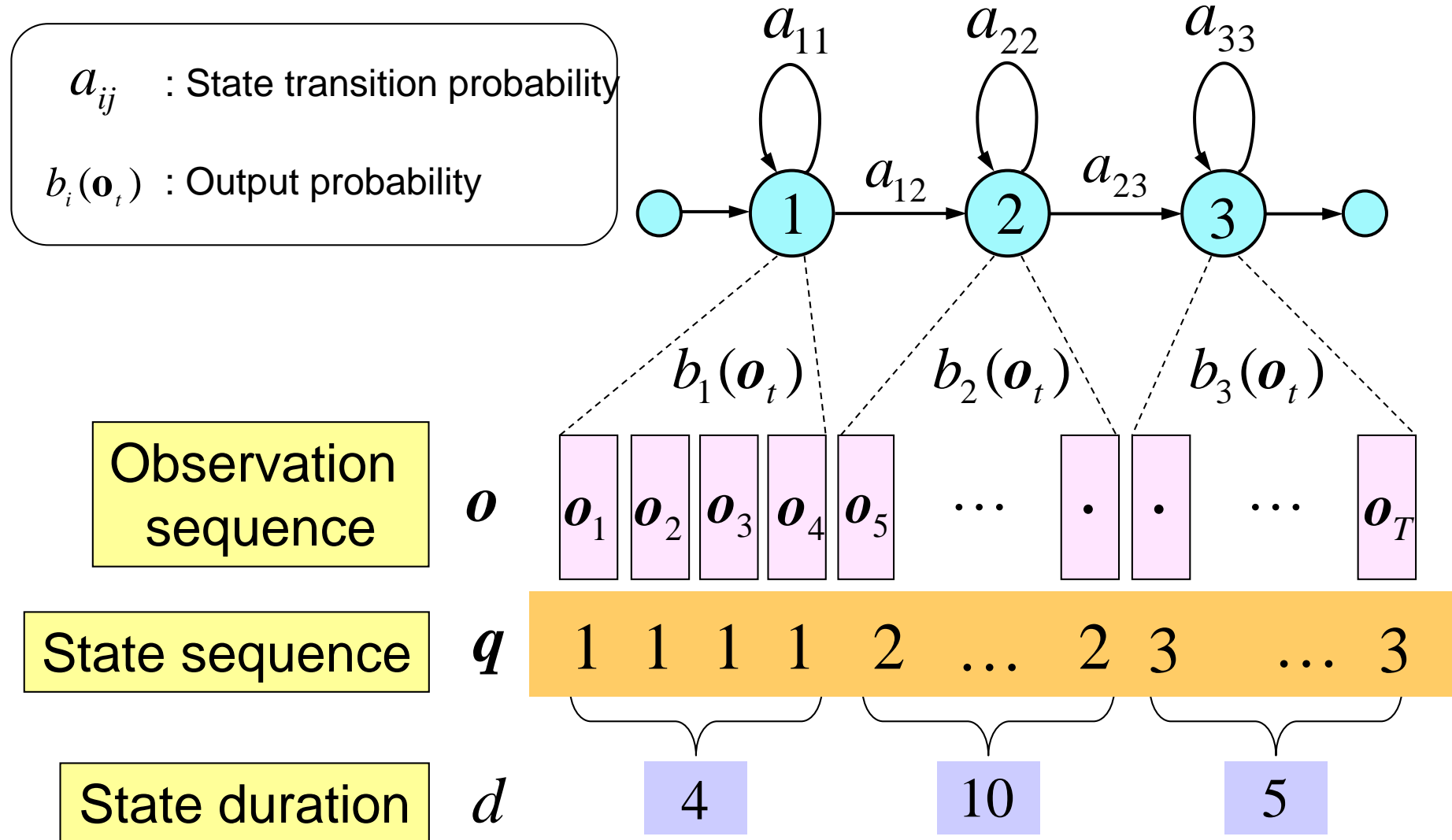
$$\begin{aligned}\hat{q} &= \arg \max_q P(q | w, \lambda) \\ \hat{\mathbf{o}} &= \arg \max_{\mathbf{o}} P(\mathbf{o} | \hat{q}, \lambda)\end{aligned}$$

Determination of State Sequence



- Speech vocoding: Source-filter model
- **Speech parameter modeling and generation with HMM**
 - Overview of HMM framework
 - State duration modeling 
 - Spectrum modeling
 - F0 modeling
 - Context clustering
- Voice character controlling
 - Adaptation (mimicking voices)
 - Interpolation (mixing voices)
- Application
 - Singing voice
 - Emotional voice
 - Audio-visual speech synthesis

Determination of State Sequence



Determine state sequence via determining state durations

Determination of State Sequence

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

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

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

d_i : state duration of i -th state

K : # of states in a sentence HMM for w

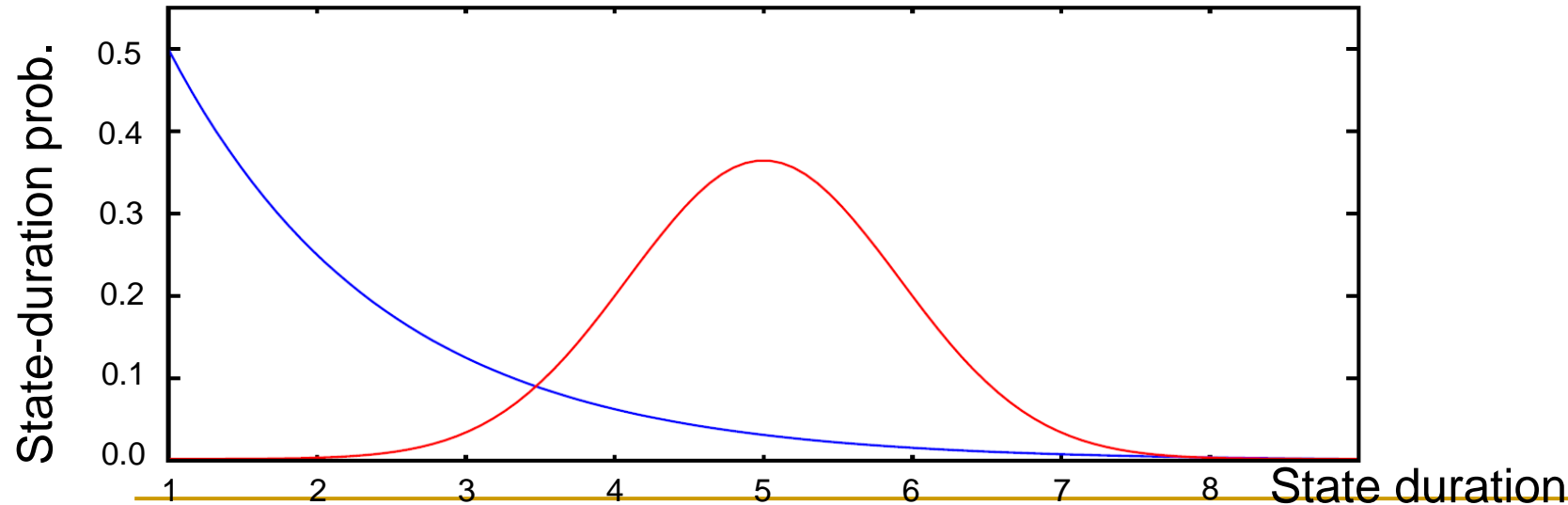
State Duration Modeling

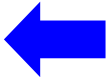
Geometric

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

Gaussian

$$p_i(d_i) = N(d_i | m_i, \sigma_i^2) \Rightarrow \hat{d}_i = m_i$$



- Speech vocoding: Source-filter model
- **Speech parameter modeling and generation with HMM**
 - Overview of HMM framework
 - Duration modeling
 - Spectrum modeling 
 - F0 modeling
 - Context clustering
- Voice character controlling
 - Adaptation (mimicking voices)
 - Interpolation (mixing voices)
- Application
 - Singing voice
 - Emotional voice
 - Audio-visual speech synthesis

Speech Parameter Generation Algorithm

For a 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

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



$$\hat{\mathbf{q}} = \arg \max_q P(\mathbf{q} | w, \lambda)$$

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

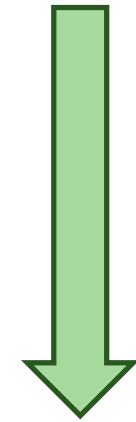
Determination of State Output

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

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

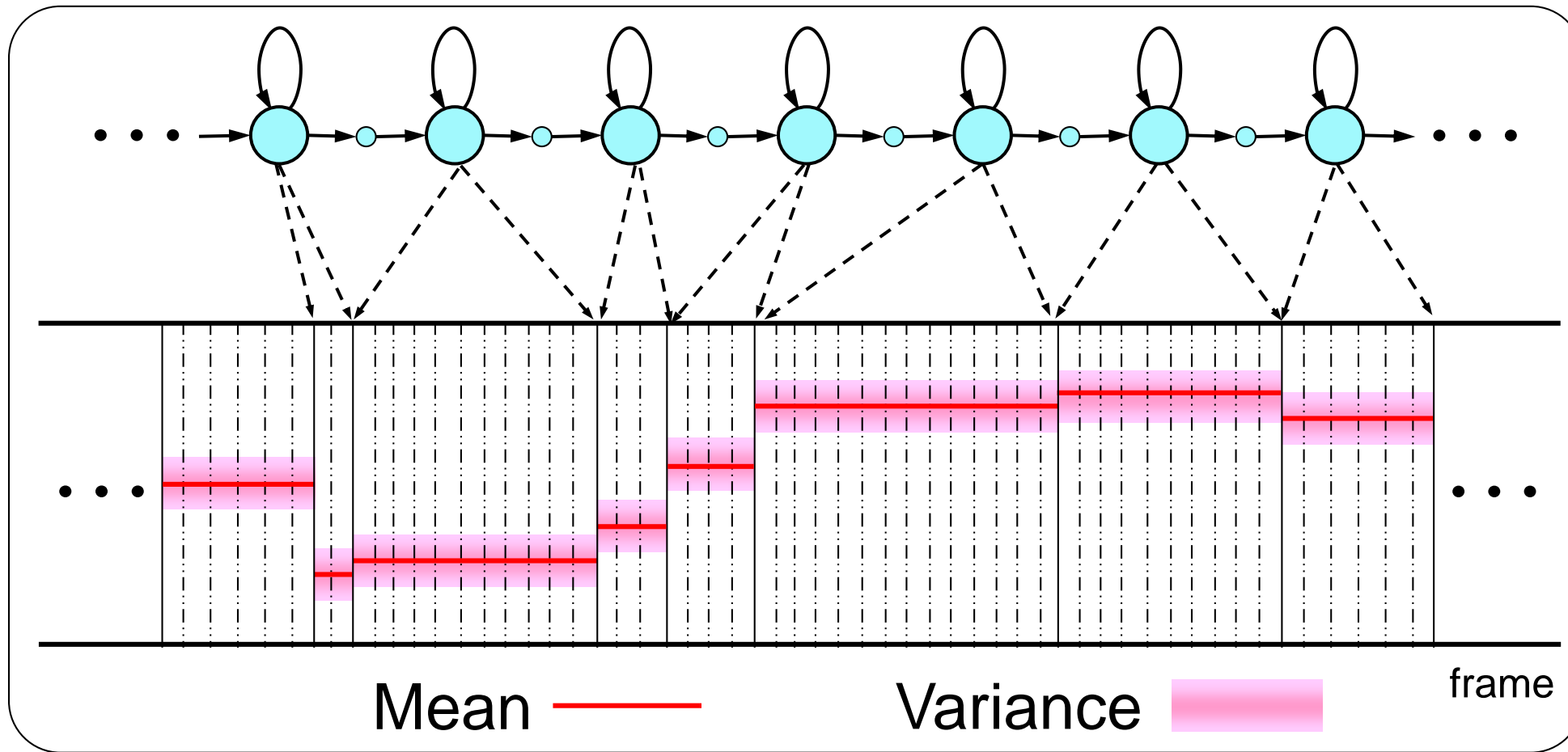
Gaussian distribution

$$b_j(\mathbf{o}_t) = N(\mathbf{o}_t | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$



$$\hat{\mathbf{o}}_t = \boldsymbol{\mu}_j$$

Generated Feature Sequence



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

To Solve the Problem ...

Constrained with dynamic features!

Let the speech parameter vector \mathbf{o}_t at frame t consists of the static feature vector \mathbf{c}_t and the dynamic feature vector $\Delta\mathbf{c}_t$:

$$\mathbf{o}_t = \{\mathbf{c}_t, \Delta\mathbf{c}_t\}$$

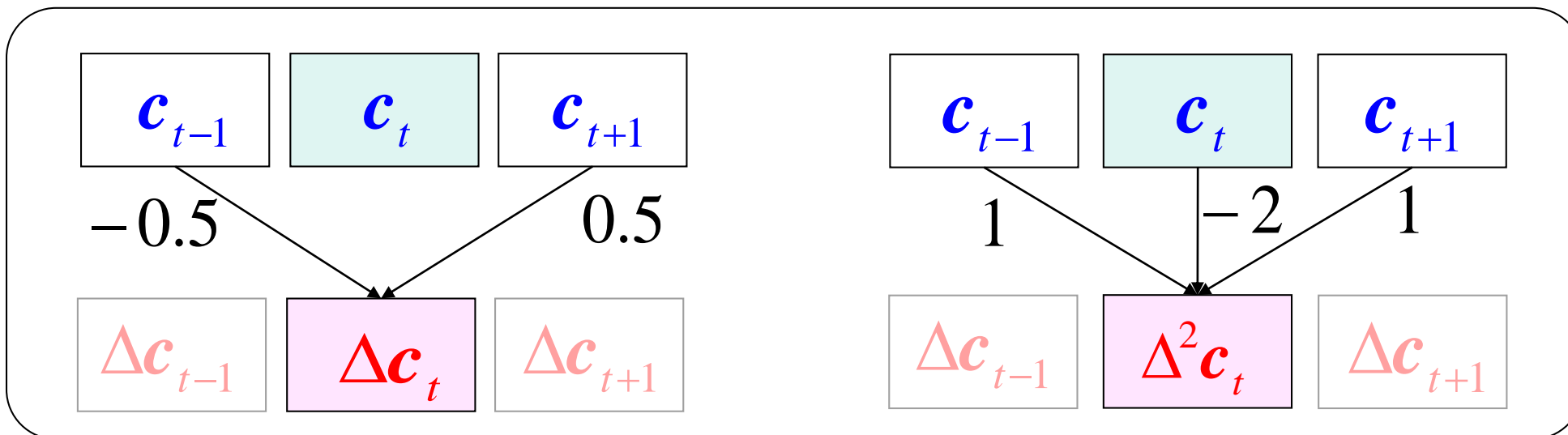
Assume \mathbf{c}_t and $\Delta\mathbf{c}_t$ are statistically independent, then

$$P(\mathbf{o} \mid \hat{\mathbf{q}}, \boldsymbol{\lambda}) = \prod_{t=1}^T b_{q_t}(\mathbf{o}_t) = \prod_{t=1}^T b_{q_t}(\mathbf{c}_t) b_{q_t}(\Delta\mathbf{c}_t)$$

Dynamic Features

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

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

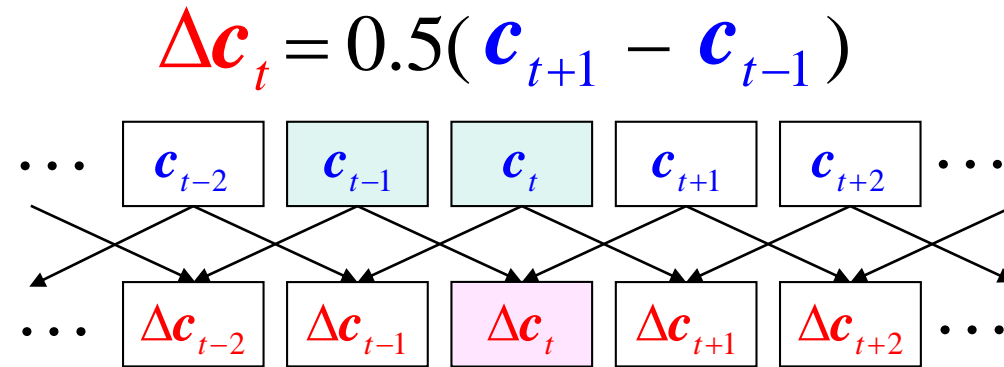


Integration of Dynamic Features

Speech parameter vector \mathbf{o}_t includes both static & dynamic features

$$\mathbf{o}_t = \begin{bmatrix} \mathbf{c}_t^\top, \Delta \mathbf{c}_t^\top \end{bmatrix}^\top$$

$\begin{bmatrix} \text{blue} \\ \text{red} \end{bmatrix} 2M$
 $\begin{bmatrix} \text{blue} \end{bmatrix} M$
 $\begin{bmatrix} \text{red} \end{bmatrix} M$

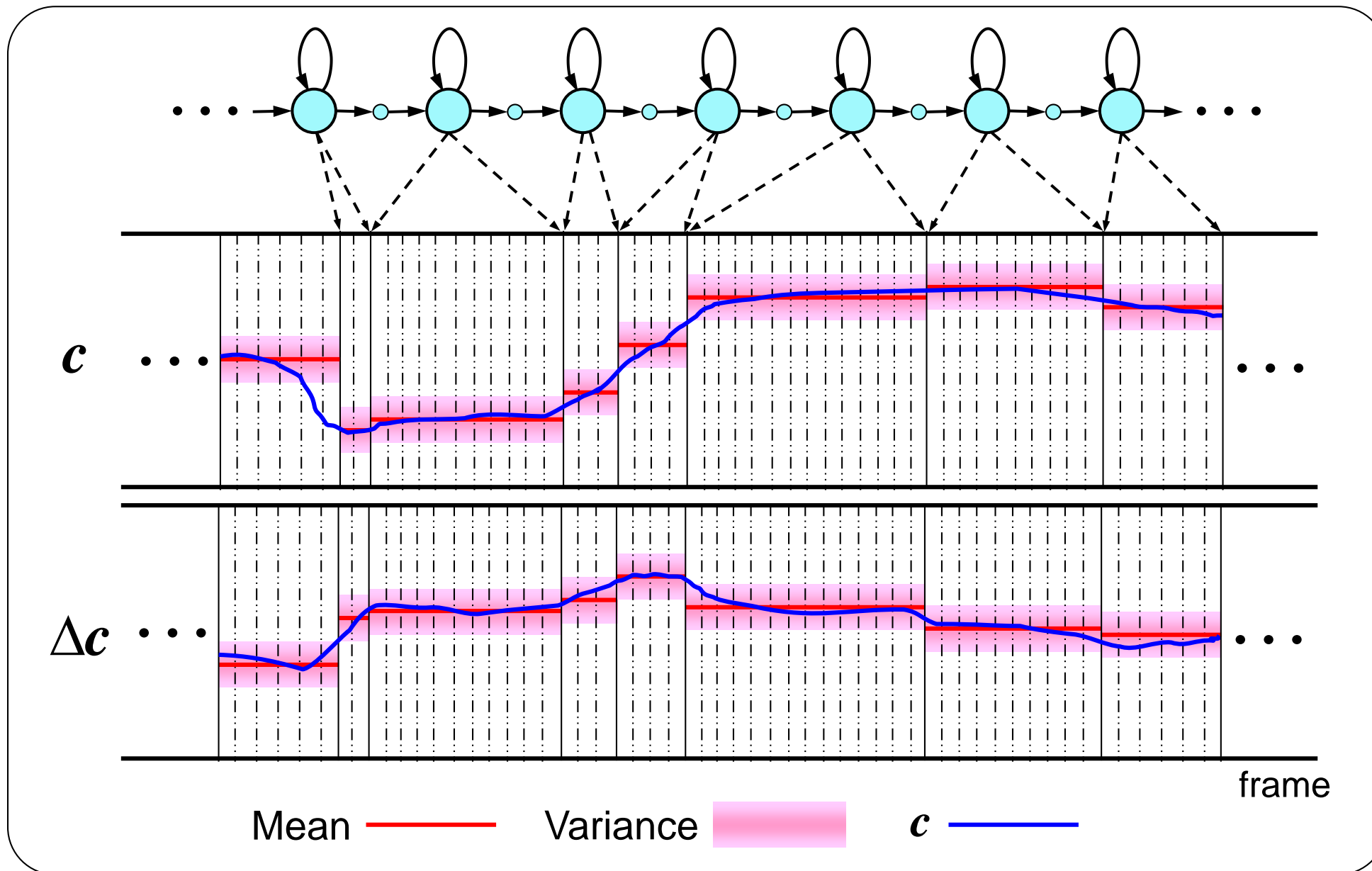


The relationship between \mathbf{o}_t & \mathbf{c}_t can be arranged as

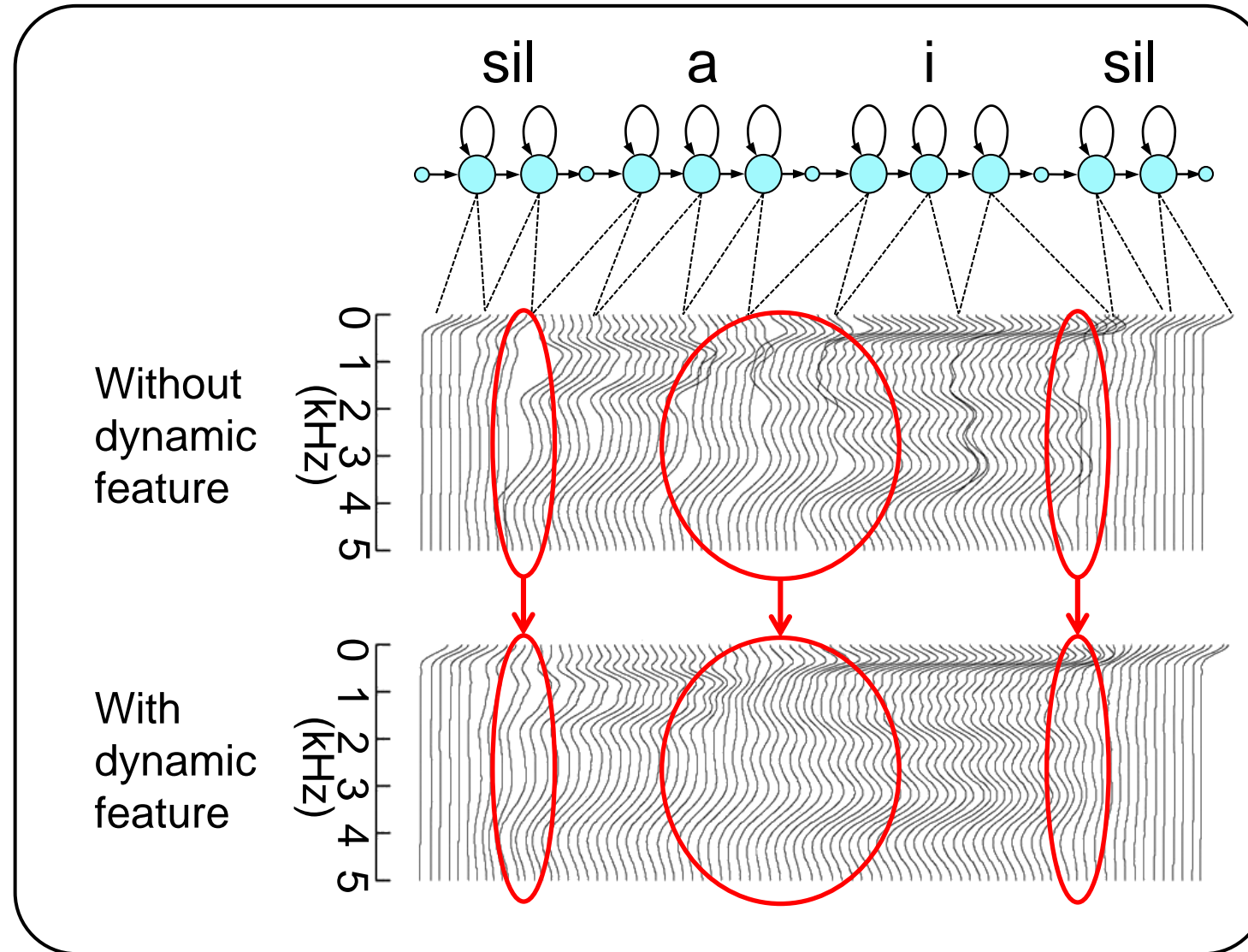
$$\begin{bmatrix} \vdots \\ \mathbf{o}_{t-1} \\ \vdots \\ \mathbf{o}_t \\ \vdots \\ \mathbf{o}_{t+1} \\ \vdots \end{bmatrix} = \begin{bmatrix} \dots & \vdots & \vdots & \vdots & \vdots & \dots \\ \dots & \mathbf{0} & \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots \\ \dots & -1/2 \mathbf{I} & \mathbf{0} & 1/2 \mathbf{I} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & -1/2 \mathbf{I} & \mathbf{0} & 1/2 \mathbf{I} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & -1/2 \mathbf{I} & \mathbf{0} & \dots \\ \dots & \vdots & \vdots & \vdots & \vdots & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ \mathbf{c}_{t-2} \\ \mathbf{c}_{t-1} \\ \mathbf{c}_t \\ \mathbf{c}_{t+1} \\ \vdots \end{bmatrix}$$

$W (M \times N)$ N

Generated Speech Parameter Trajectory

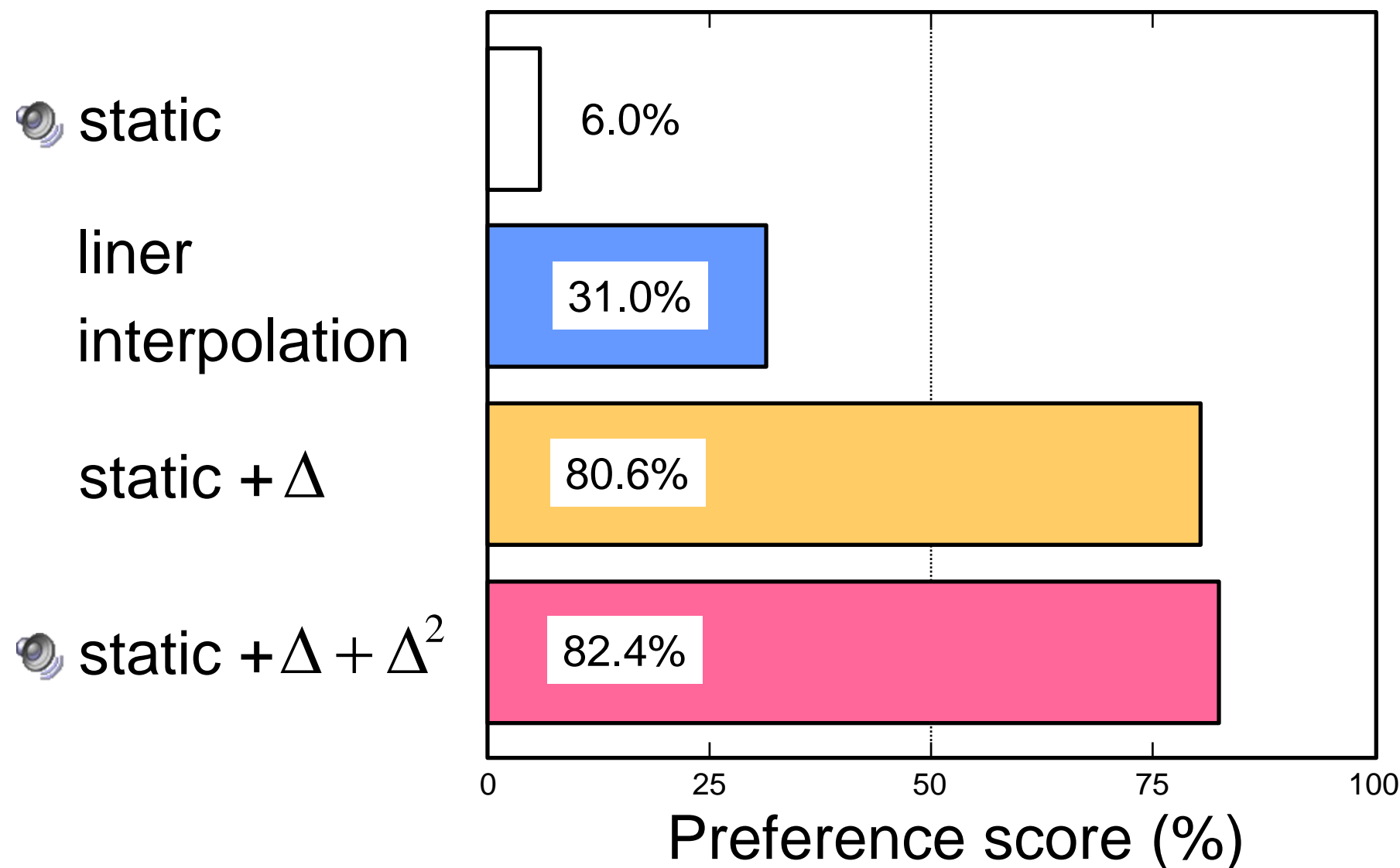


Generated Spectra



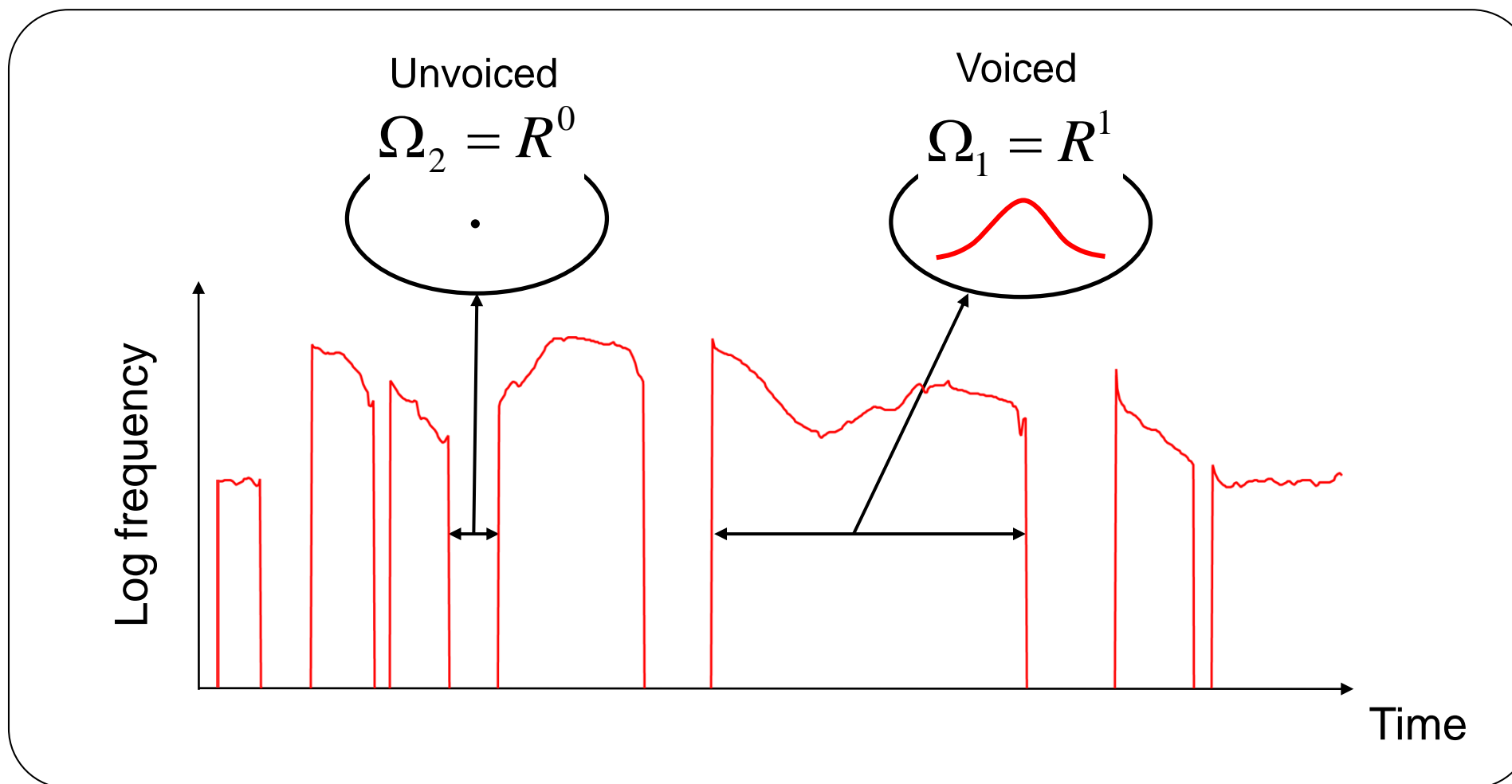
Spectra changing smoothly between phonemes

Effect of Dynamic Features (Japanese)



- Speech vocoding: Source-filter model
- **Speech parameter modeling and generation with HMM**
 - Overview of HMM framework
 - Duration modeling
 - Spectrum modeling
 - F0 modeling ←
 - Context clustering
- Voice character controlling
 - Adaptation (mimicking voices)
 - Interpolation (mixing voices)
- Application
 - Singing voice
 - Emotional voice
 - Audio-visual speech synthesis

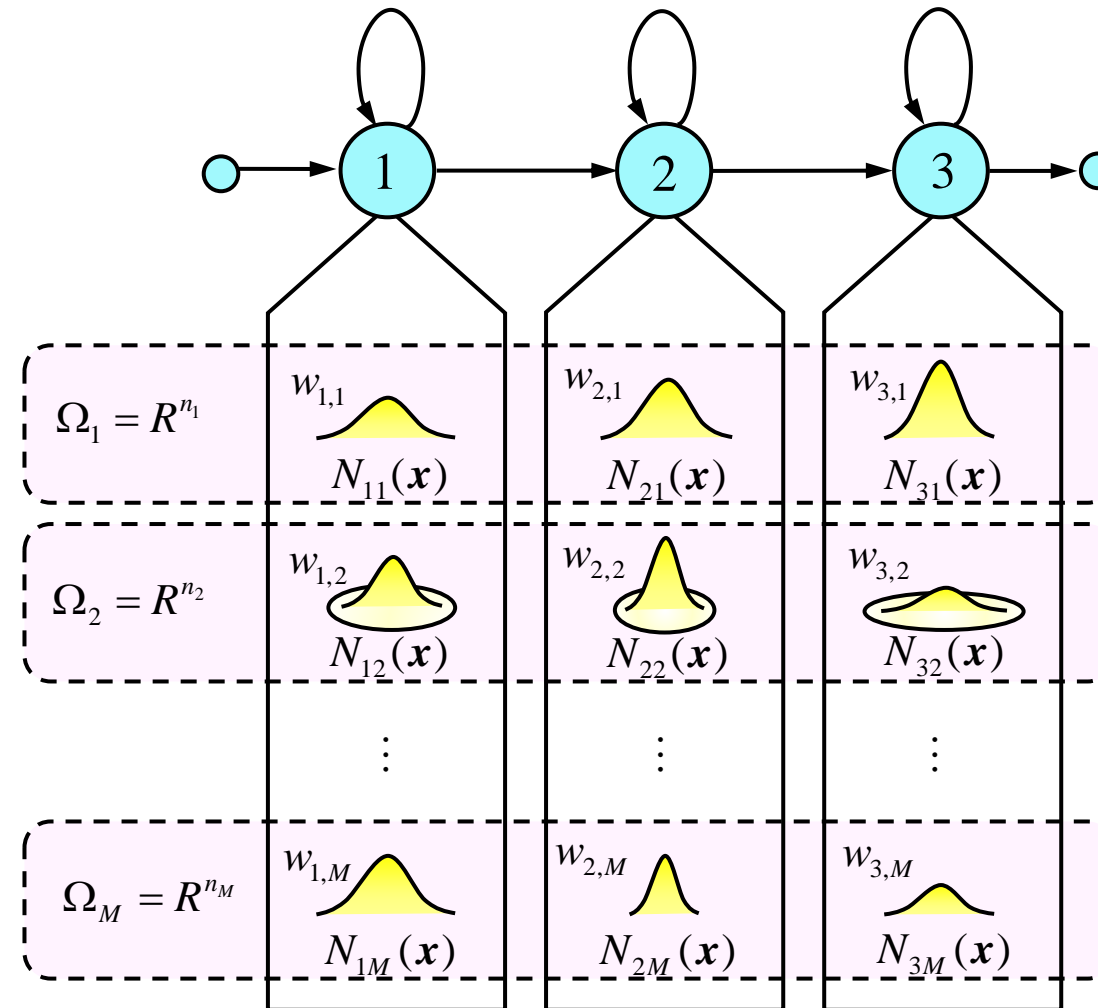
Observation of F0



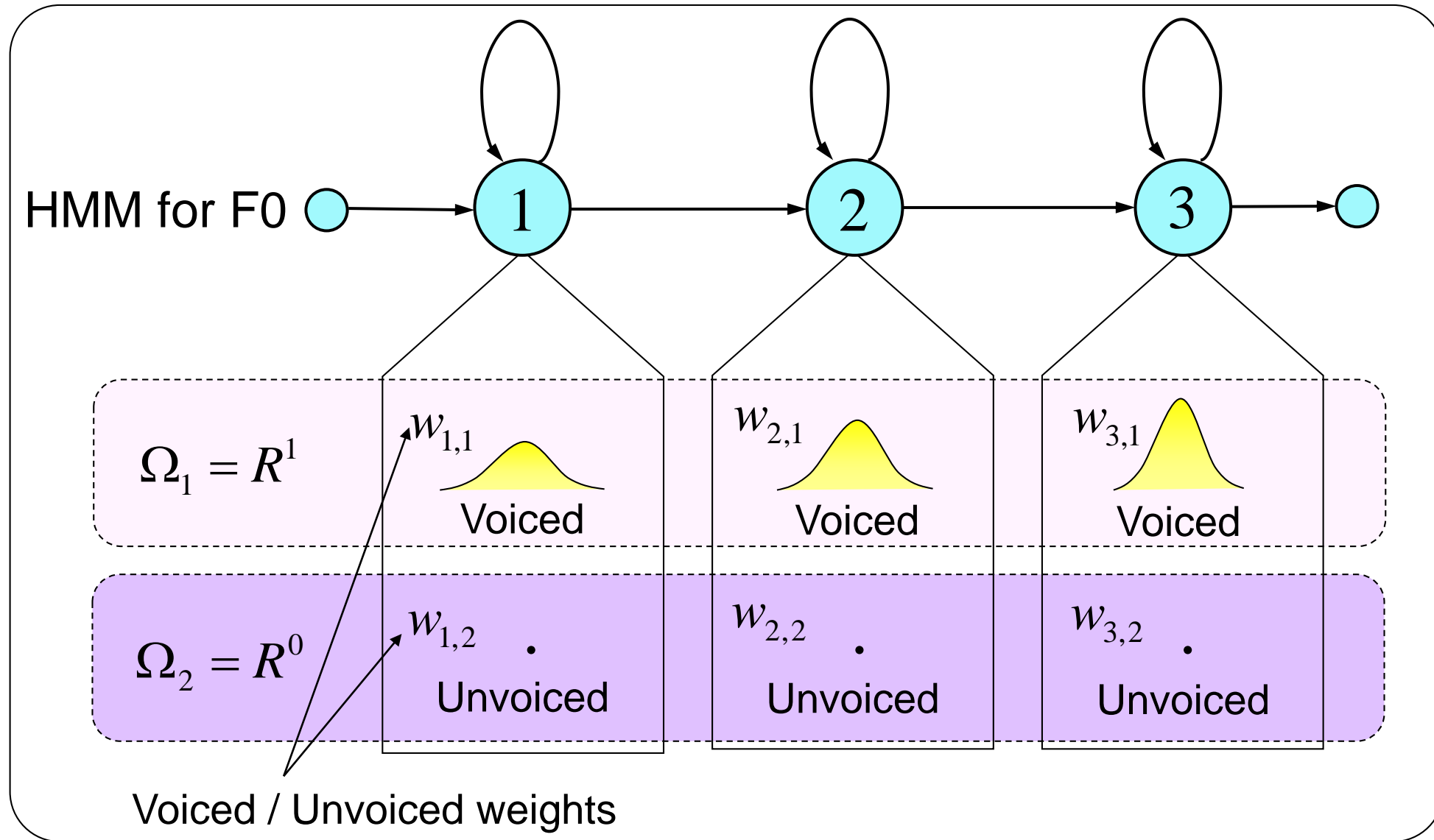
Unable to model by continuous or discrete distribution
⇒ Multi-space probability distribution HMM (MSD-HMM)

Structure of MSD-HMM

Each state has two or more distributions

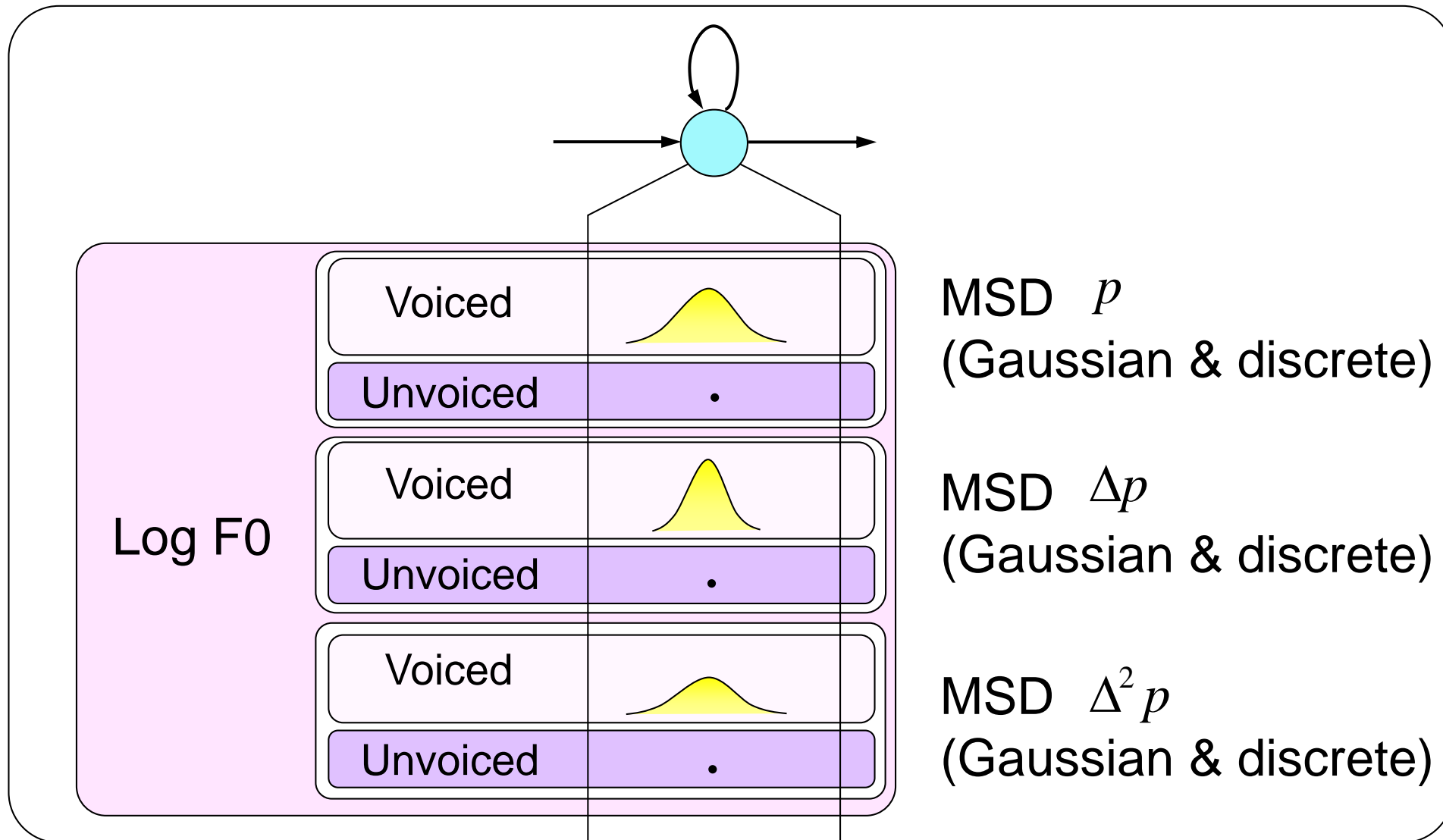


MSD-HMM for F0 Modeling



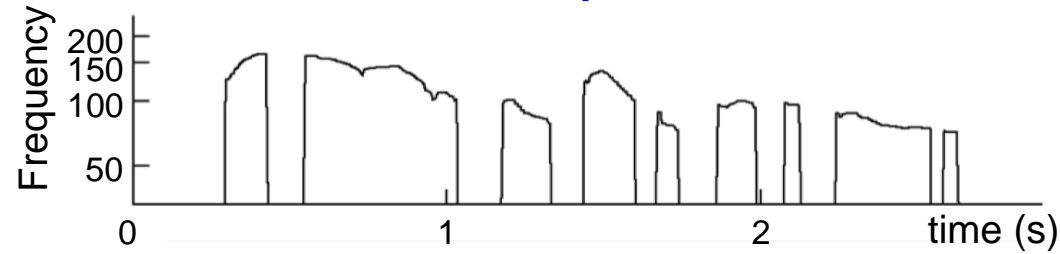
To Model Dynamic F0 Features

■ Structure of F0 state-output distributions

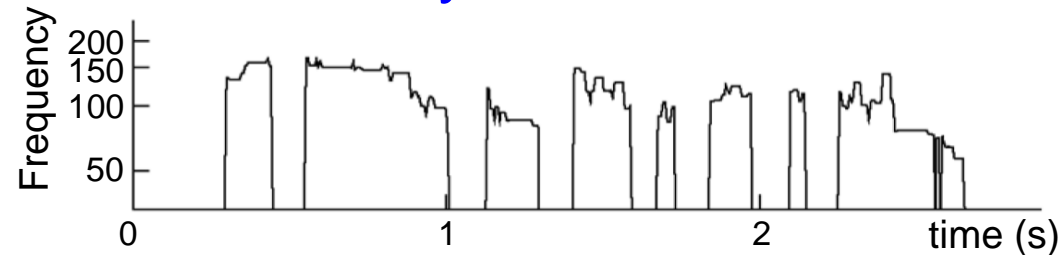


Generated F0 Trajectory

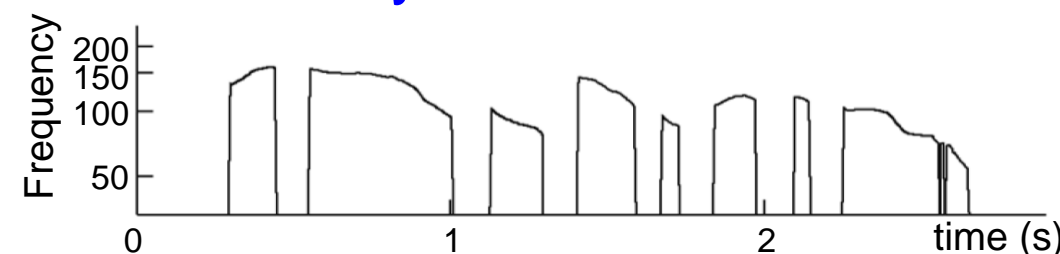
natural speech



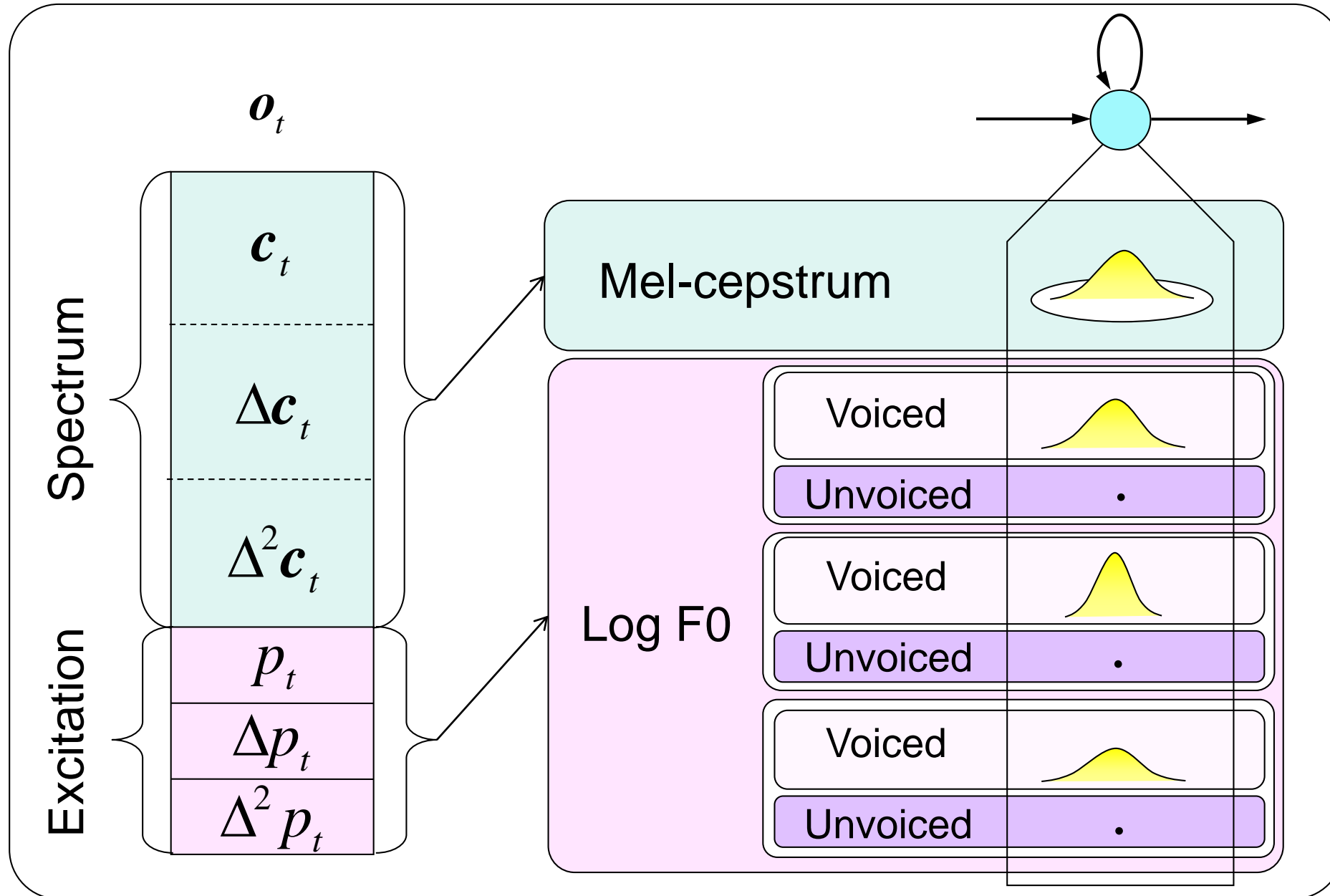
without dynamic features












with dynamic features

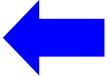


To Summarize: Structure of State-Output Distribution

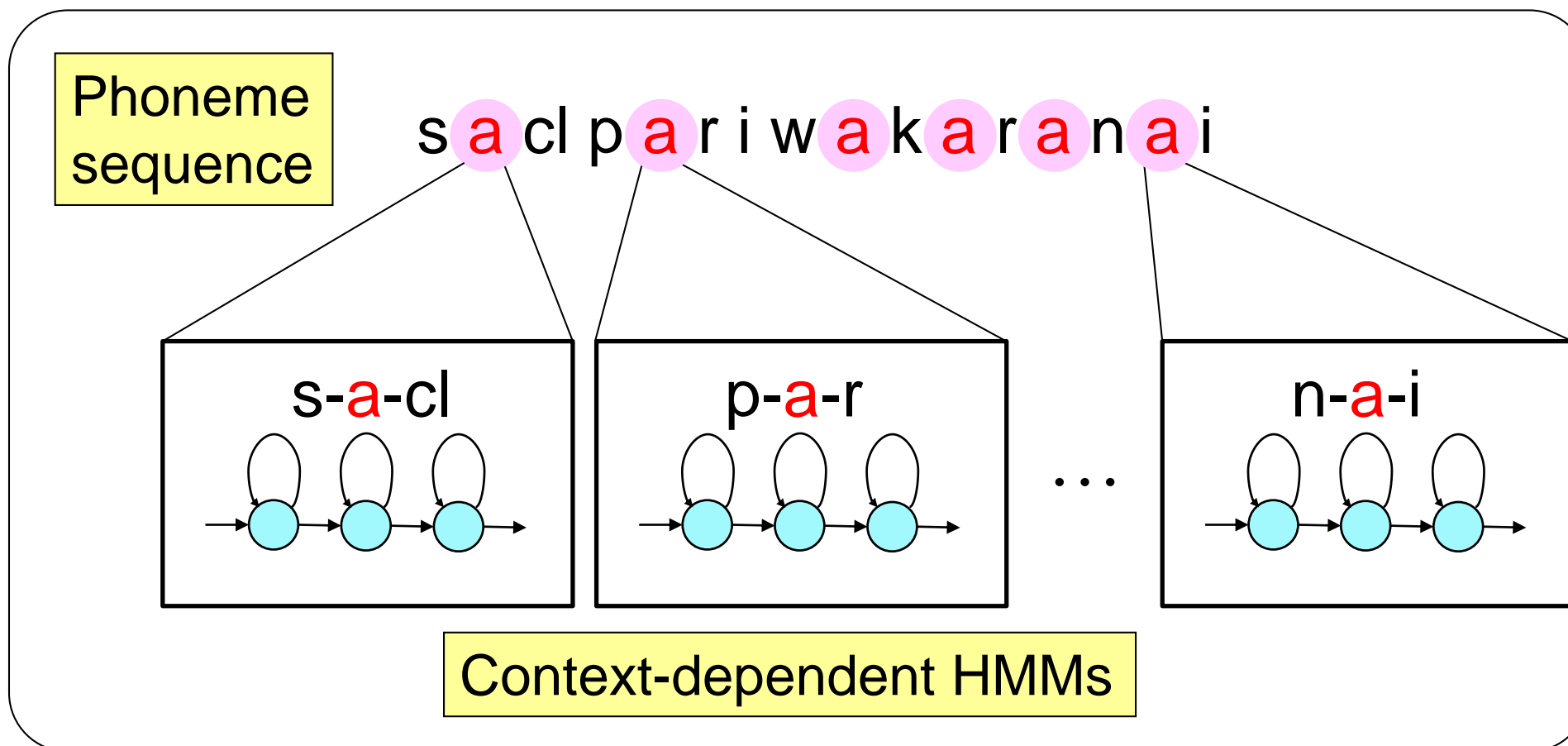


To Summarize: Speech Samples (Japanese)

		Mel-cepstrum		
		static	static + Δ	static + $\Delta + \Delta^2$
log F0	static			
	static + Δ			
	static + $\Delta + \Delta^2$			

- Speech vocoding: Source-filter model
- **Speech parameter modeling and generation with HMM**
 - Overview of HMM framework
 - Duration modeling
 - Spectrum modeling
 - F0 modeling
 - Context clustering 
- Voice character controlling
 - Adaptation (mimicking voices)
 - Interpolation (mixing voices)
- Application
 - Singing voice
 - Emotional voice
 - Audio-visual speech synthesis

Context-Dependent Model



- **Considering relations between phonemes**
 - Context \Rightarrow factor of speech variations
 - Improving model accuracy

Phoneme

- {preceding, succeeding} two phonemes
- current phoneme

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

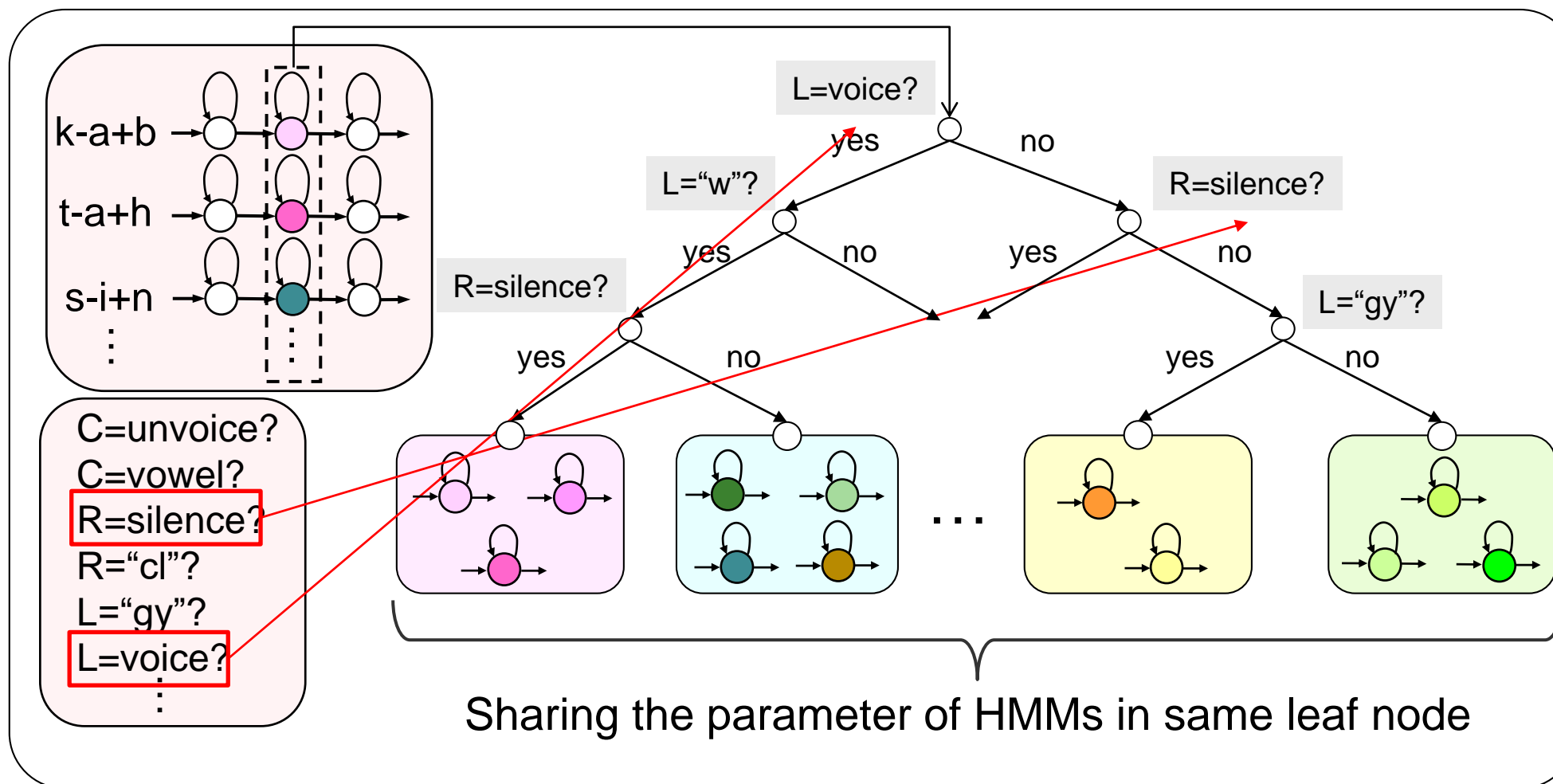
Phrase

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

.....

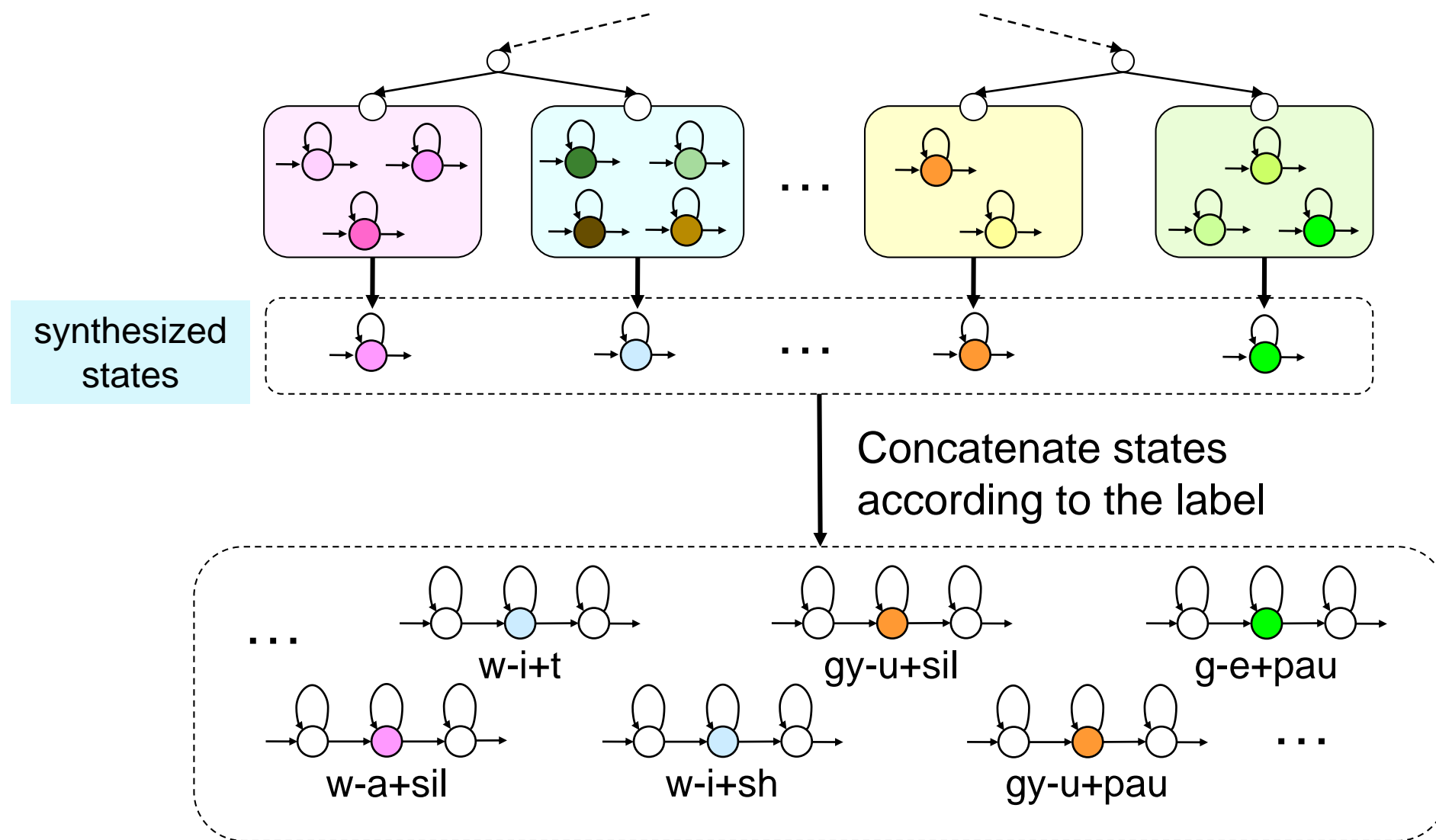
Huge # of combinations \Rightarrow Difficult to have all possible models

Decision Tree-based State Clustering [Odell; '95]

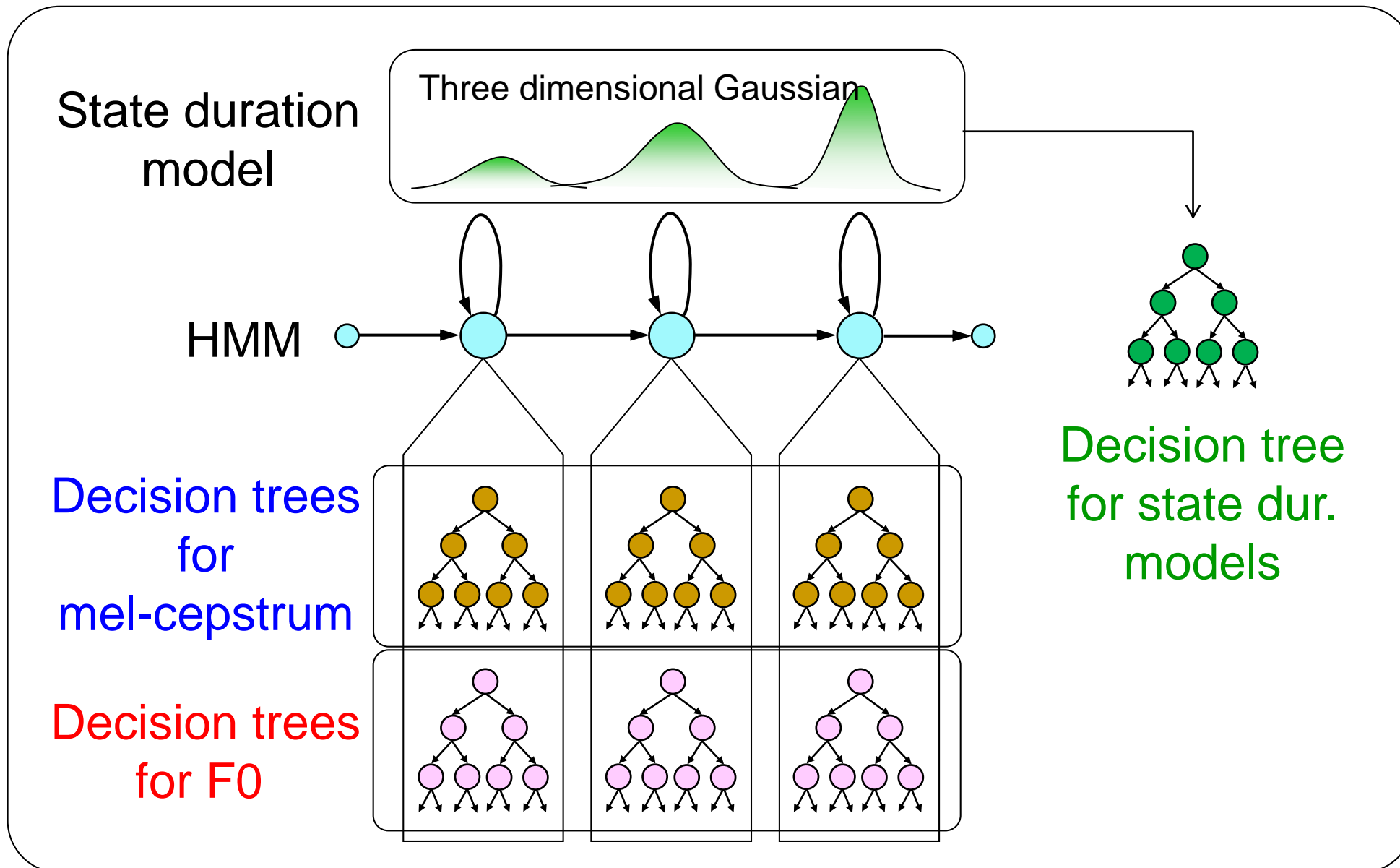


- Each state separated automatically by the optimum question
- The optimum question determined for increasing likelihood

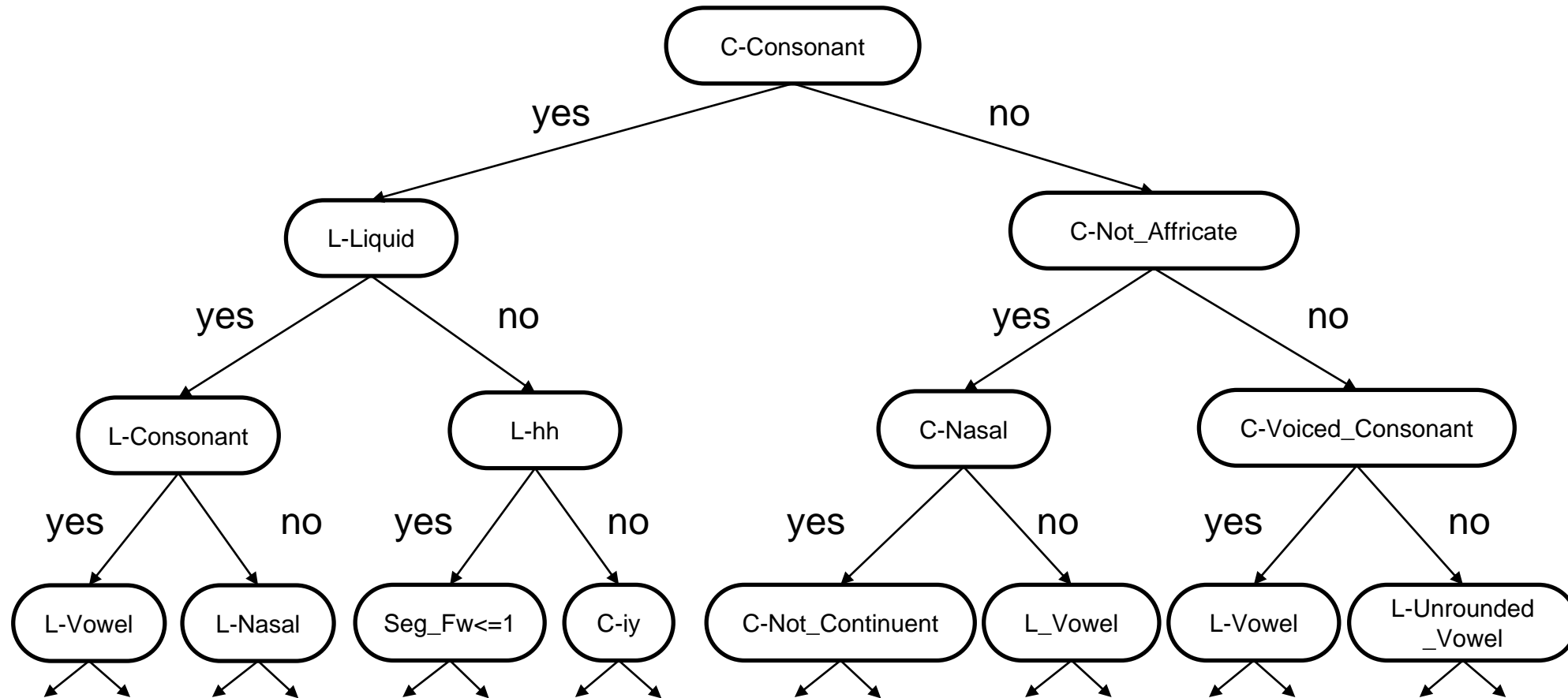
Synthesize From Leaf Nodes



Stream-dependent Tree-based Clustering

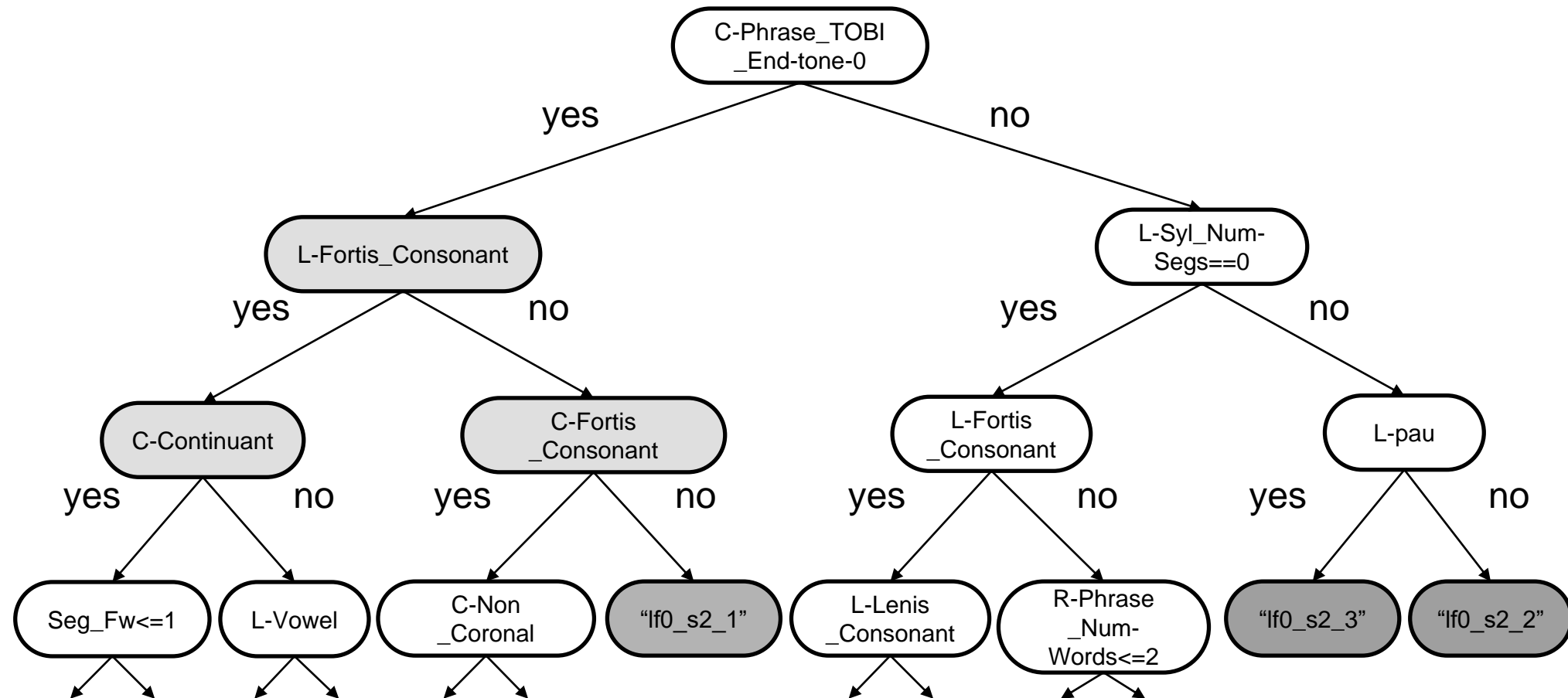


Tree for Spectrum (1st state)



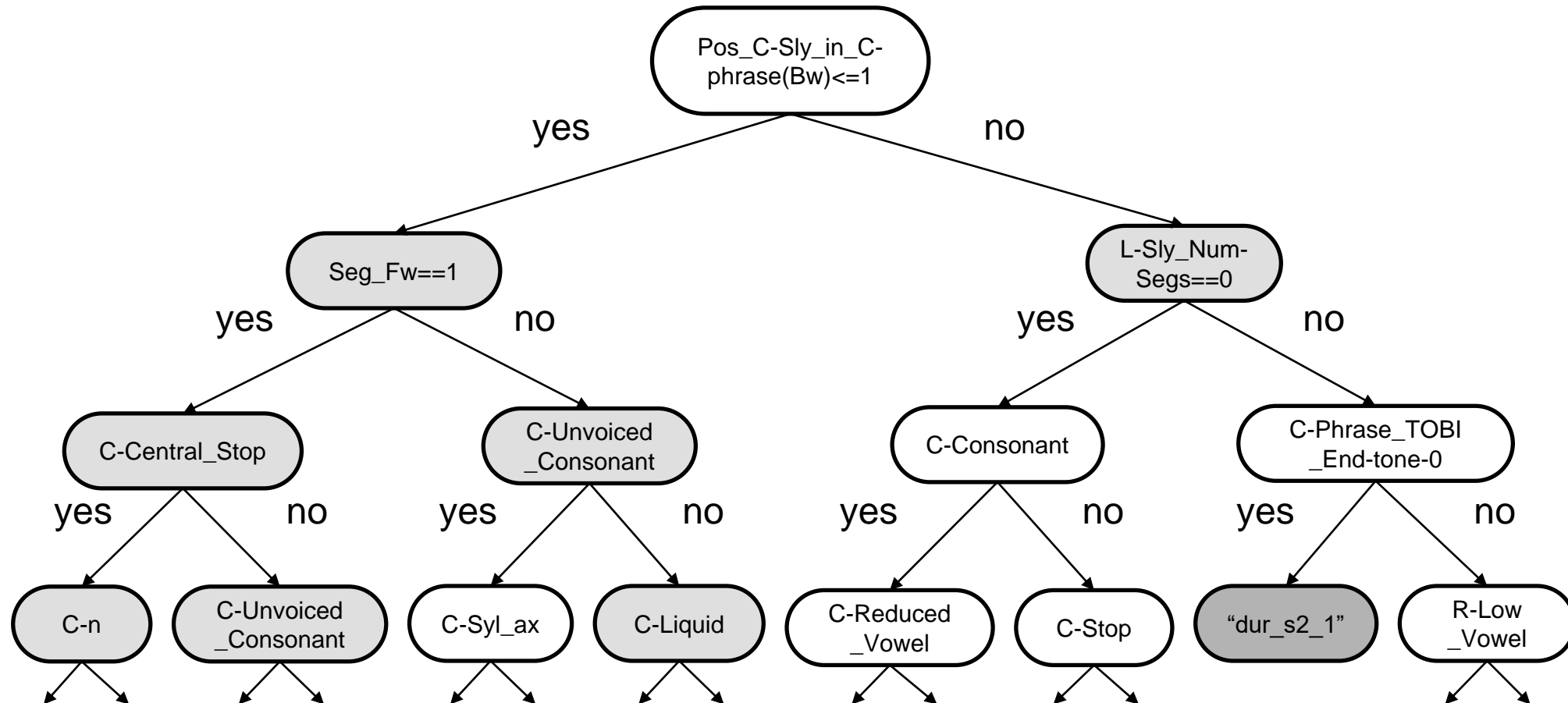
- Questions about phonetic attributes

Tree for F0 (1st state)



– Questions about linguistic attributes

Tree for State Duration



- Linguistic questions for pause
- Phonetic questions for speech

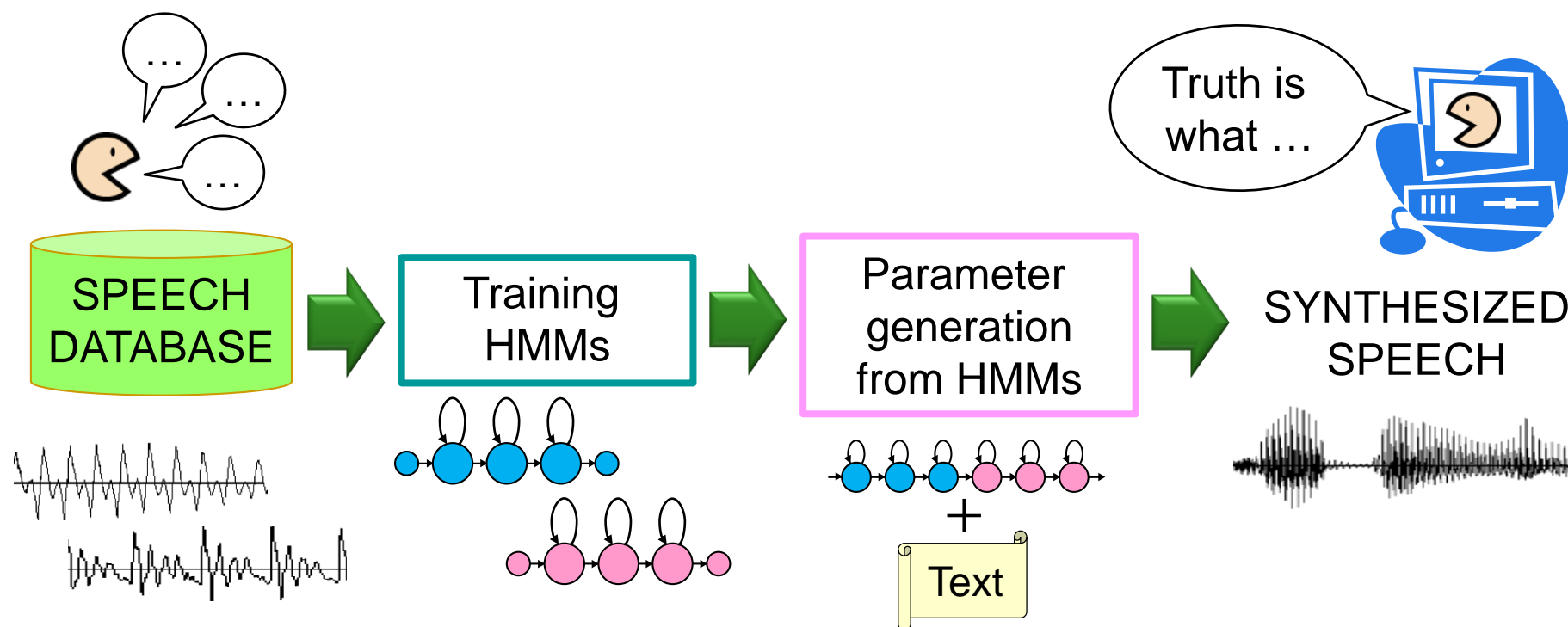
Speech Synthesis from HMM

- **Speech parameters are generated from HMMs**

- Spectrum parameters
- Excitation parameters (F0)

- **Vocoding parameters to synthesize speech**

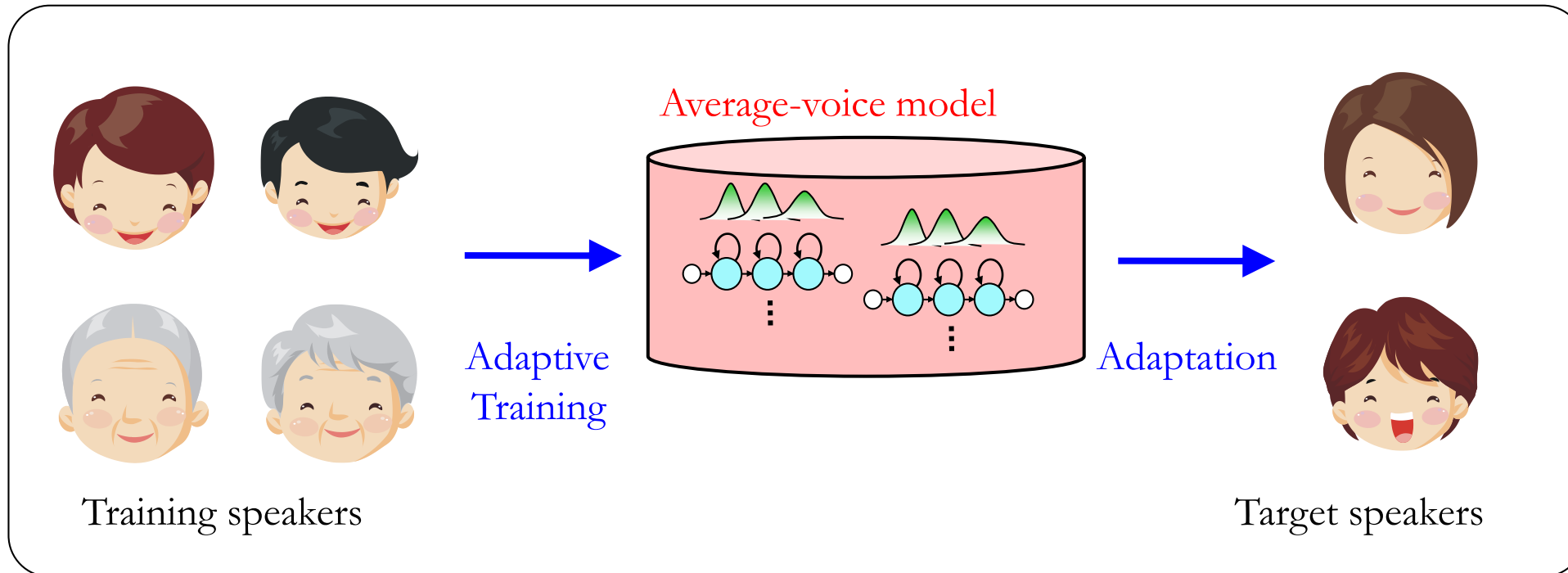
⇒ Obtain high-quality synthesized speech



- **Speech vocoding: Source-filter model**
- **Speech parameter modeling and generation with HMM**
 - ❑ Overview of HMM framework
 - ❑ State duration modeling
 - ❑ Spectrum modeling
 - ❑ F0 modeling
 - ❑ Context clustering
- **Voice character controlling**
 - ❑ Adaptation (mimicking voices)
 - ❑ Interpolation (mixing voices)
- **Application**
 - ❑ Singing voice
 - ❑ Emotional voice
 - ❑ Audio-visual speech synthesis

Adaptation (Mimicking Voices)

- **Adaptation / adaptive training of HMMs**
 - Originally developed in ASR, but works very well in TTS
 - Average voice-based speech synthesis (AVSS) [Yamagishi; '06]








- Require small data of target speaker / speaking style
⇒ small cost to create new voices

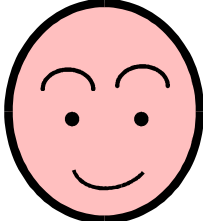
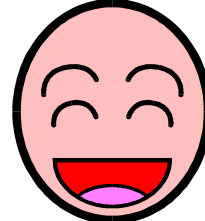

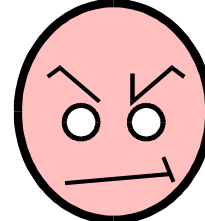







Adaptation Demo

■ Speaker adaptation

Original voice: 

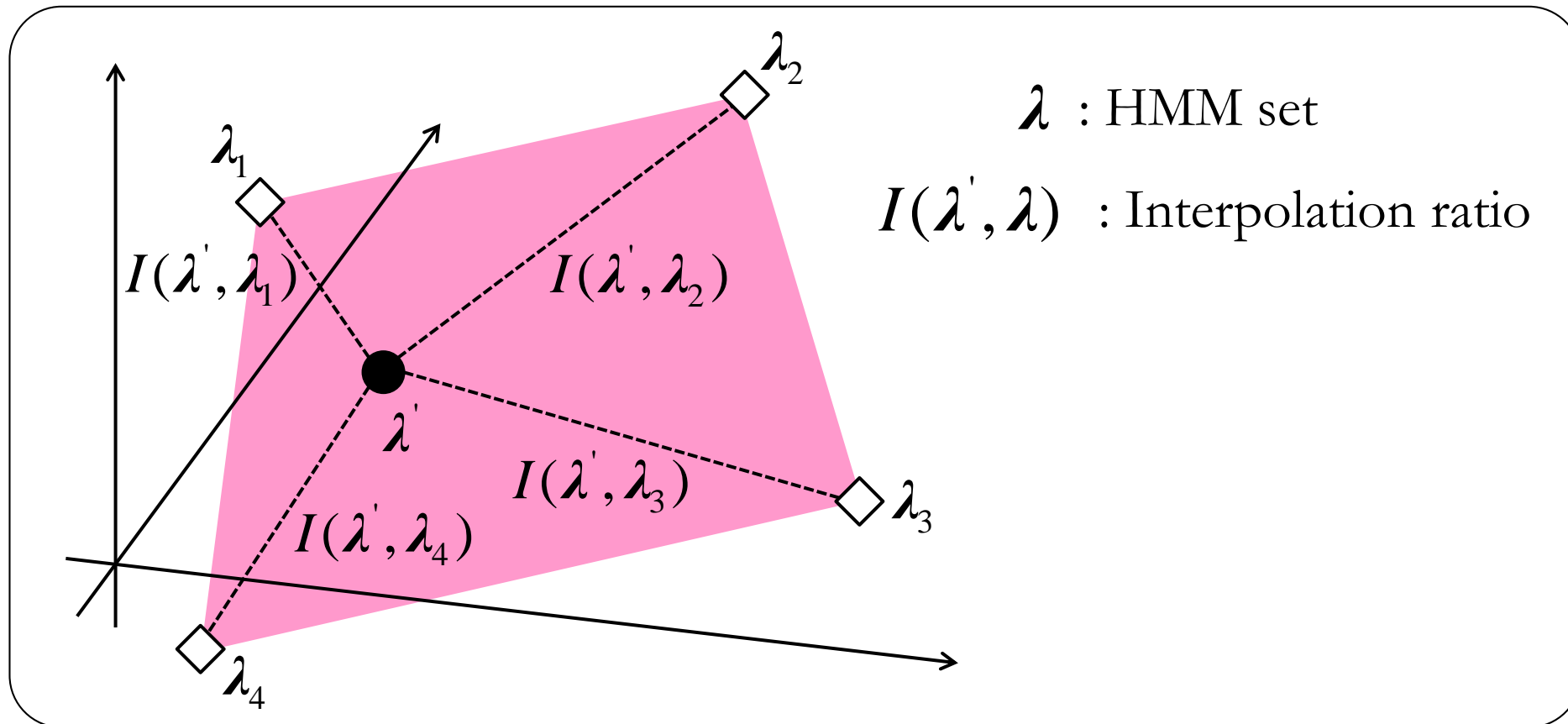
Average voice model	Number of adaptation sentences			
	10 sentences	100 sentences	500 sentences	1132 sentences
				

■ Style adaptation

				
	Neutral	Joyful	Sad	Rough
English				
Japanese				

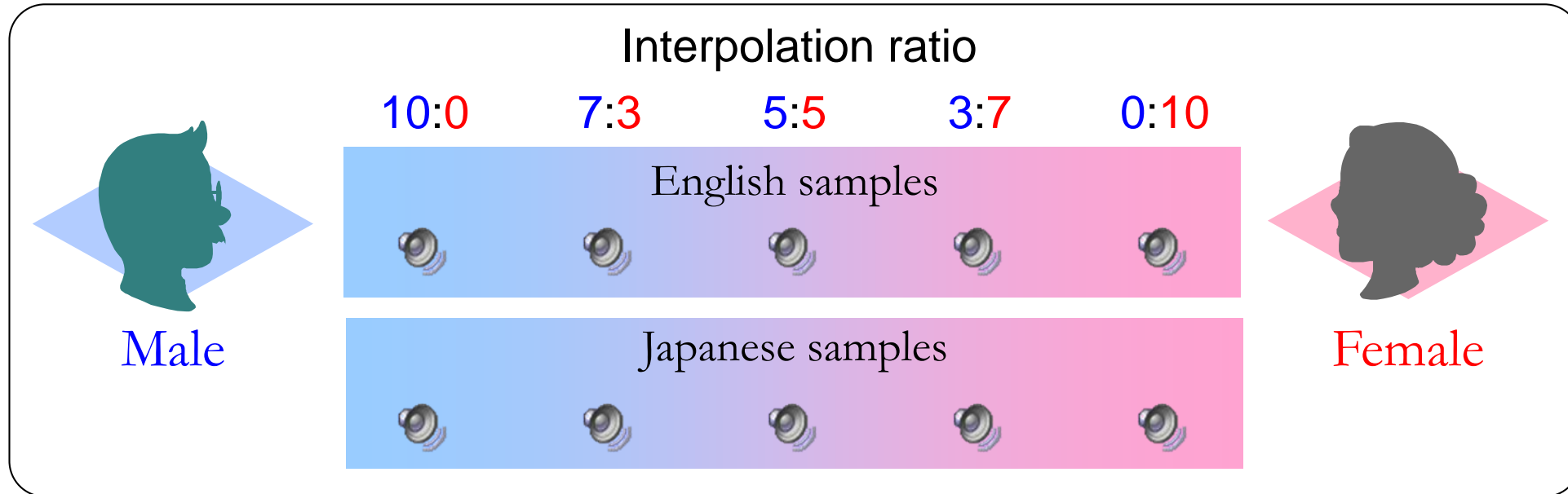
Interpolation (Mixing Voices)

- **Interpolate parameters among representative HMM sets**
 - Create new voices even if no adaptation data is available
 - Gradually change speaker & speaking styles [Yoshimura; '97, Tachibana; '05]

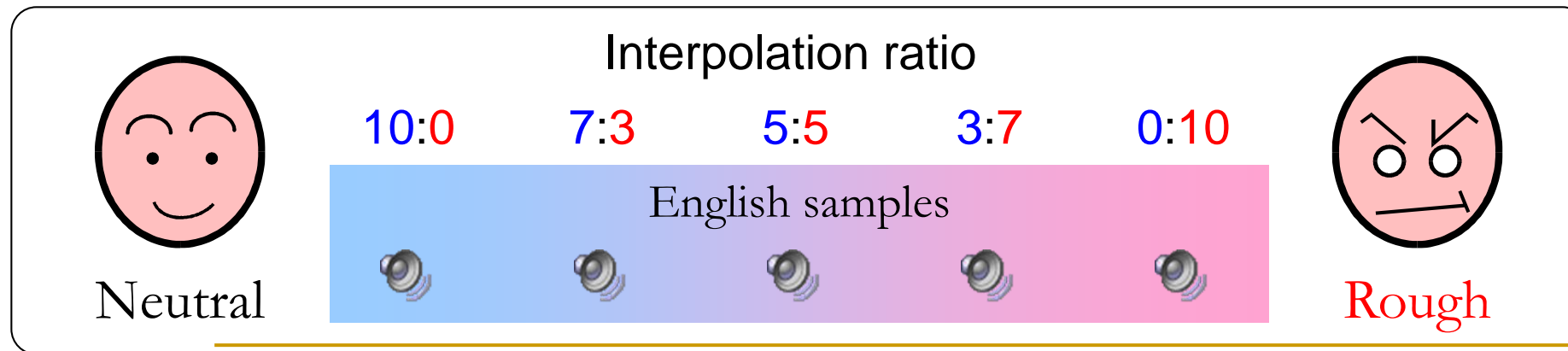


Interpolation Demo

■ Speaker interpolation

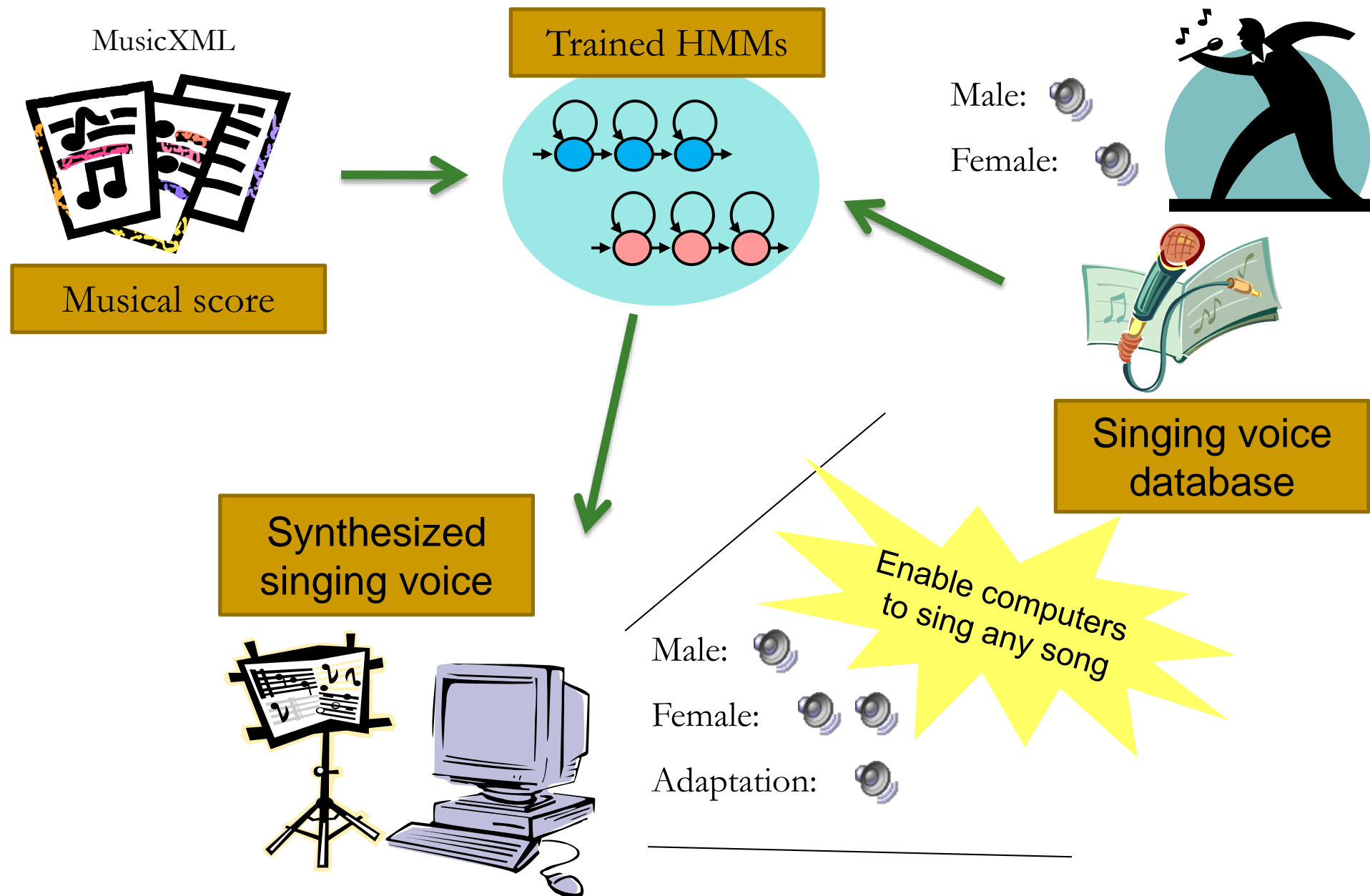


■ Style interpolation

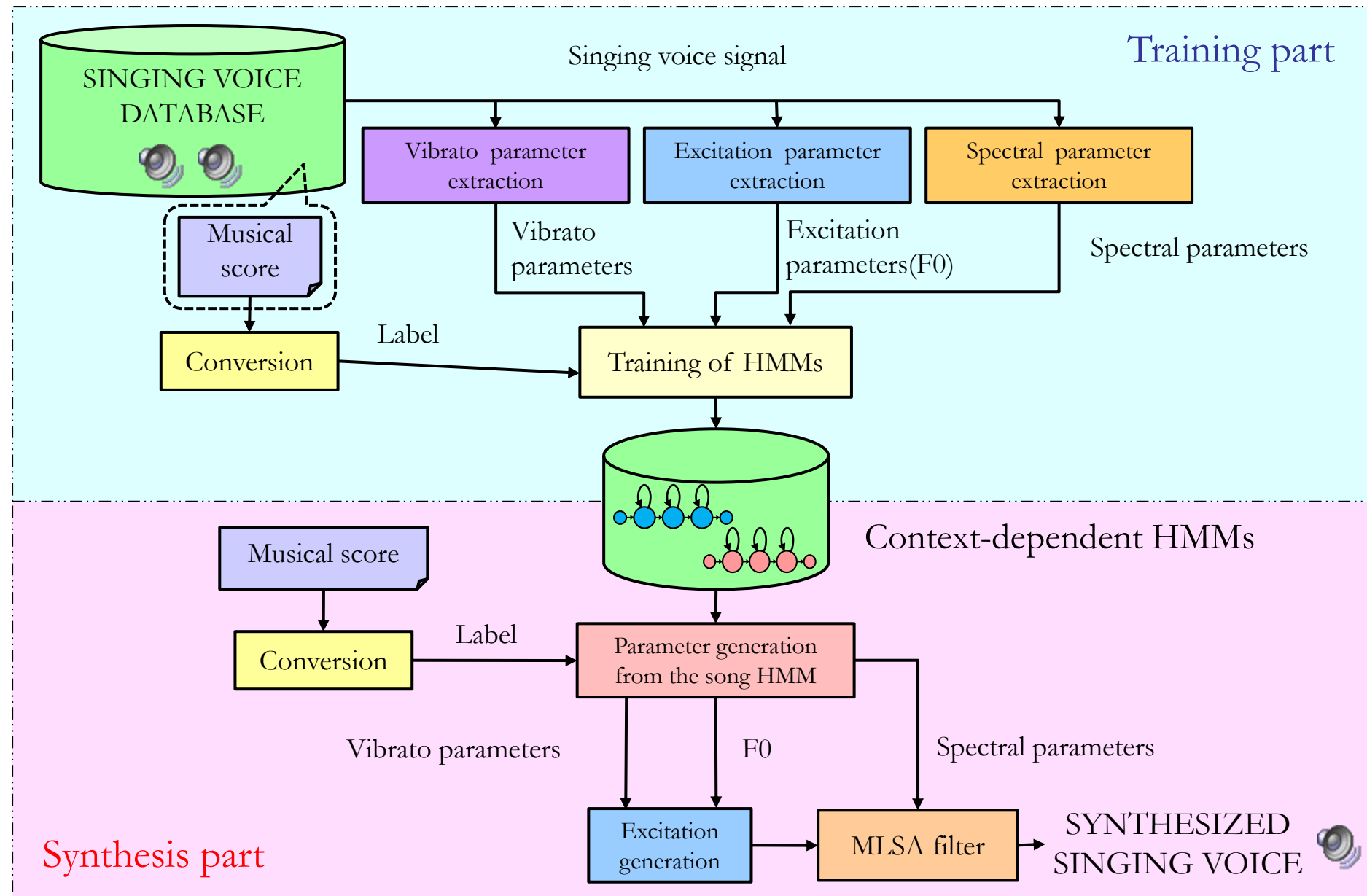


- **Speech vocoding: Source-filter model**
- **Speech parameter modeling and generation with HMM**
 - ❑ Overview of HMM framework
 - ❑ State duration modeling
 - ❑ Spectrum modeling
 - ❑ F0 modeling
 - ❑ Context clustering
- **Voice character controlling**
 - ❑ Adaptation (mimicking voices)
 - ❑ Interpolation (mixing voices)
- **Application**
 - ❑ Singing voice
 - ❑ Emotional voice
 - ❑ Audio-visual speech synthesis

Singing Voice Synthesis [Oura; '10]

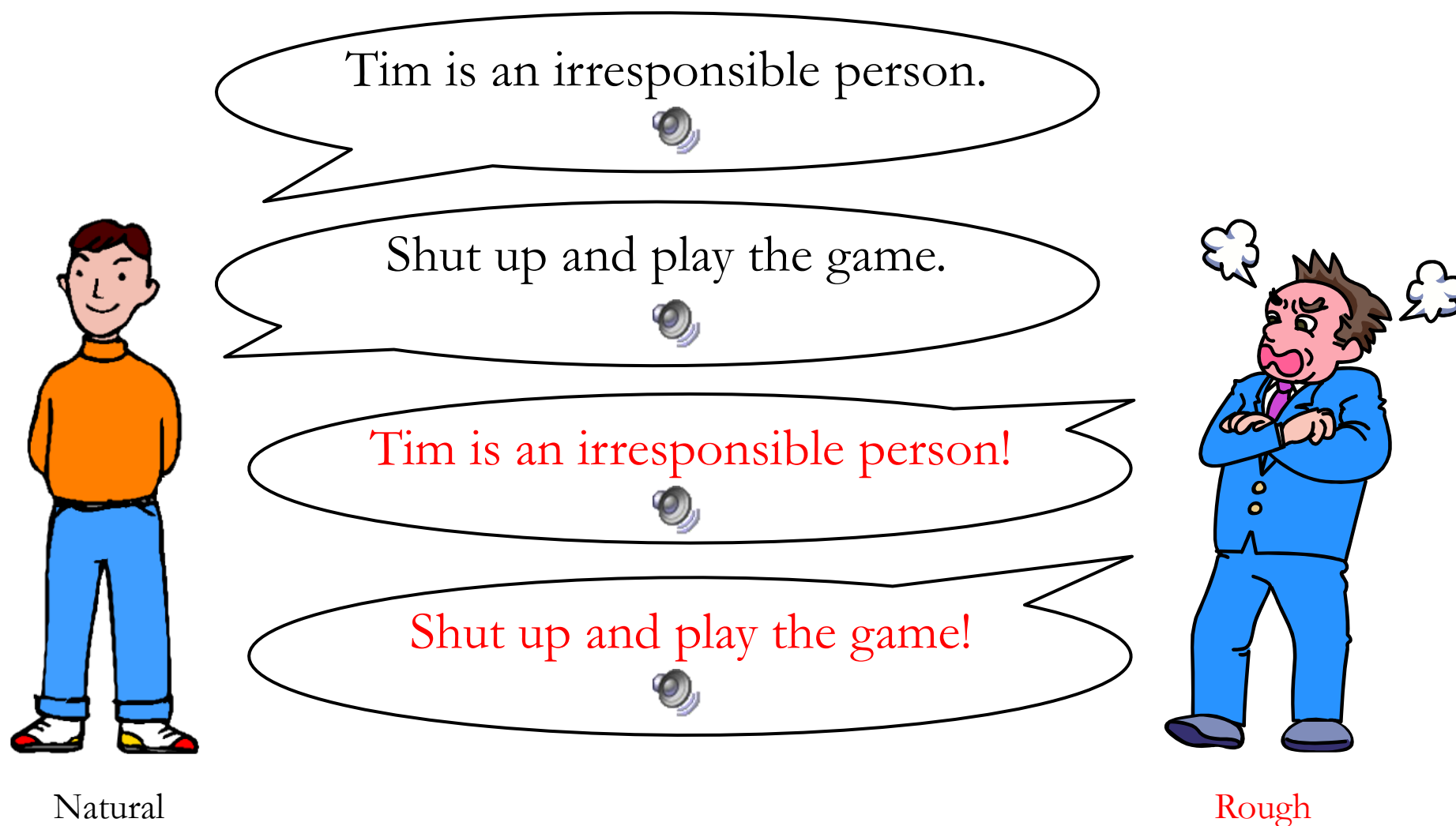


Singing Voice Synthesis



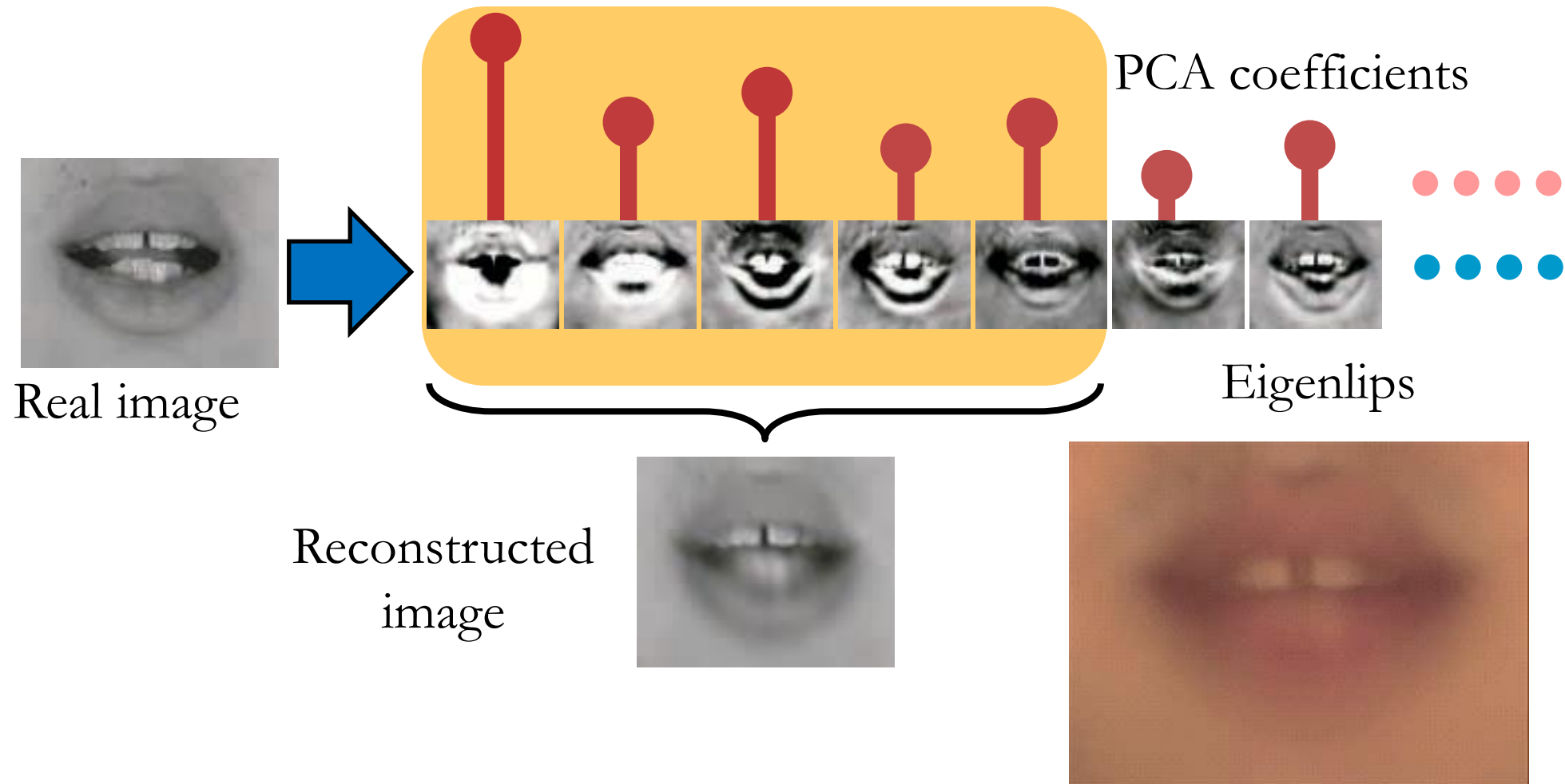
Emotional Speech Synthesis [Tsuzuki; '04]

Using emotional voices for training data

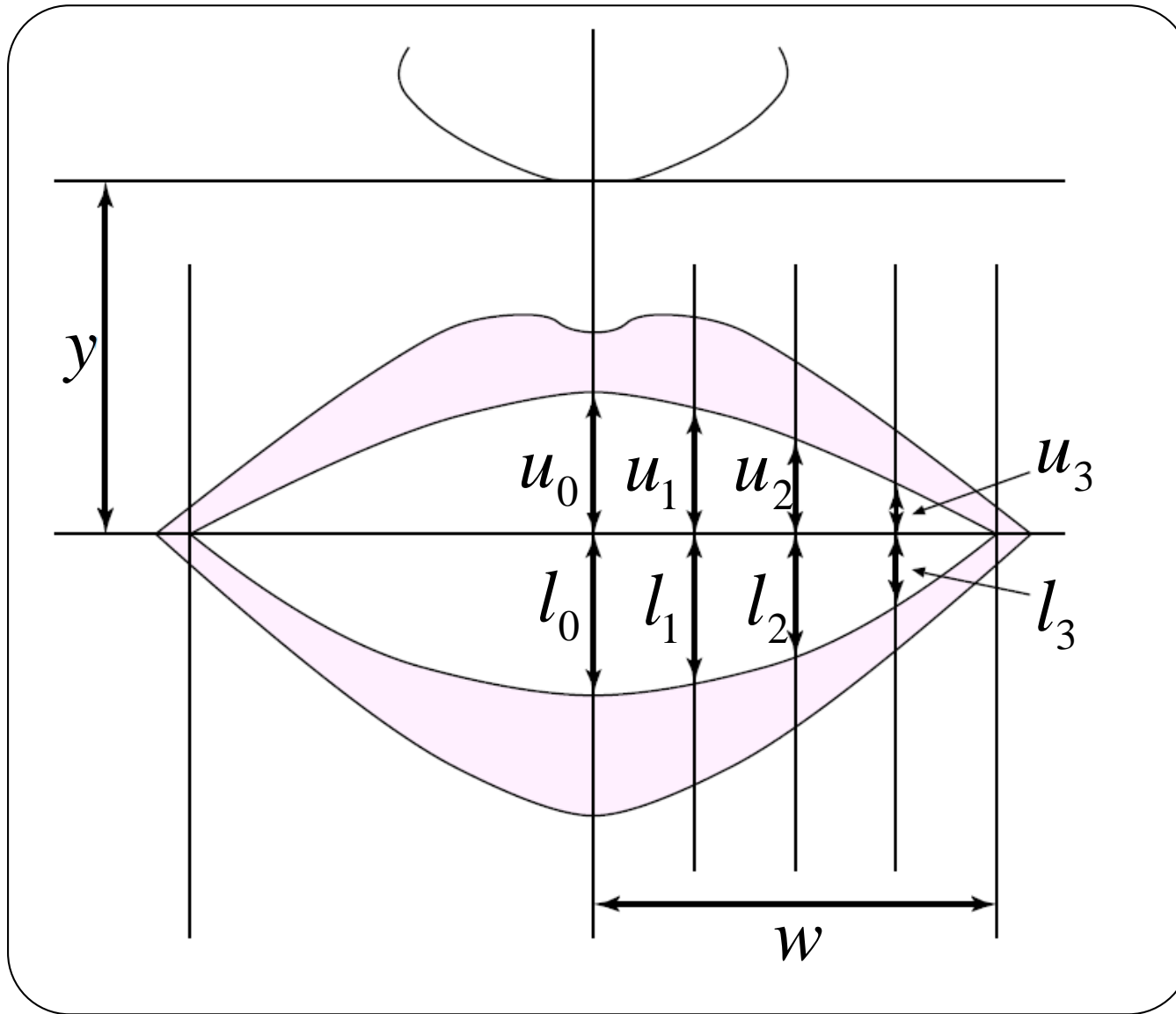


Audio-Visual Synthesis (Pixel-based) [Sako; '00]

- Pixel image: high dimensionality
⇒ Dimensionality reduction by PCA



Audio-Visual Synthesis (Model-based) [Tamura; '98]



References (1/4)

- Sagisaka;'92 - "ATR nu-TALK speech synthesis system," ICSLP, '92.
- Black;'96 - "Automatically clustering similar units...", Euro speech, '97.
- Beutnagel;'99 - "The AT&T Next-Gen TTS system," Joint ASA, EAA, & DAEA meeting, '99.
- Yoshimura;'99 - "Simultaneous modeling of spectrum ...," Eurospeech, '99.
- Itakura;'70 - "A statistical method for estimation of speech spectral density...", Trans. IEICE, J53-A, '70.
- Imai;'88 - "Unbiased estimator of log spectrum and its application to speech signal...", EURASIP, '88.
- Kobayashi;'84 - "Spectral analysis using generalized cepstrum," IEEE Trans. ASSP, 32, '84.
- Tokuda;'94 - "Mel-generalized cepstral analysis -- A unified approach to speech spectral...", ICSLP, '94.
- Imai;'83 - "Cepstral analysis synthesis on the mel frequency scale," ICASSP, '83.
- Fukada;'92 - "An adaptive algorithm for mel-cepstral analysis of speech," ICASSP, '92.
- Itakura;'75 - "Line spectrum representation of linear predictive coefficients of speech...", J. ASA (57), '75.
- Tokuda;'02 - "Multi-space probability distribution HMM," IEICE Trans. E85-D(3), '02.
- Odell;'95 - "The use of context in large vocabulary...", PhD thesis, University of Cambridge, '95.
- Shinoda;'00 - "MDL-based context-dependent subword modeling...", Journal of ASJ(E) 21(2), '00.
- Yoshimura;'98 - "Duration modeling for HMM-based speech synthesis," ICSLP, '98.
- Tokuda;'00 - "Speech parameter generation algorithms for HMM-based speech synthesis," ICASSP, '00.
- Kobayashi;'85 - "Mel generalized-log spectrum approximation...", IEICE Trans. J68-A (6), '85.
- Hunt;'96 - "Unit selection in a concatenative speech synthesis system using...", ICASSP, '96.
- Donovan;'95 - "Improvements in an HMM-based speech synthesiser," Eurospeech, '95.
- Kawai;'04 - "XIMERA: A new TTS from ATR based on corpus-based technologies," ISCA SSW5, '04.
- Hirai;'04 - "Using 5 ms segments in concatenative speech synthesis," Proc. ISCA SSW5, '04.

References (2/4)

- Rouibia;'05 - "Unit selection for speech synthesis based on a new acoustic target cost," Interspeech, '05.
- Huang;'96 - "Whistler: A trainable text-to-speech system," ICSLP, '96.
- Mizutani;'02 - "Concatenative speech synthesis based on HMM," ASJ autumn meeting, '02.
- Ling;'07 - "The USTC and iFlytek speech synthesis systems...", Blizzard Challenge workshop, 07.
- Ling;'08 - "Minimum unit selection error training for HMM-based unit selection...", ICASSP, 08.
- Plumpe;'98 - "HMM-based smoothing for concatenative speech synthesis," ICSLP, '98.
- Wouters;'00 - "Unit fusion for concatenative speech synthesis," ICSLP, '00.
- Okubo;'06 - "Hybrid voice conversion of unit selection and generation...", IEICE Trans. E89-D(11), '06.
- Aylett;'08 - "Combining statistical parametric speech synthesis and unit selection..." LangTech, '08.
- Pollet;'08 - "Synthesis by generation and concatenation of multiform segments," Interspeech, '08.
- Yamagishi;'06 - "Average-voice-based speech synthesis," PhD thesis, Tokyo Inst. of Tech., '06.
- Yoshimura;'97 - "Speaker interpolation in HMM-based speech synthesis system," Eurospeech, '97.
- Tachibana;'05 - "Speech synthesis with various emotional expressions...", IEICE Trans. E88-D(11), '05.
- Kuhn;'00 - "Rapid speaker adaptation in eigenvoice space," IEEE Trans. SAP 8(6), '00.
- Shichiri;'02 - "Eigenvoices for HMM-based speech synthesis," ICSLP, '02.
- Fujinaga;'01 - "Multiple-regression hidden Markov model," ICASSP, '01.
- Nose;'07 - "A style control technique for HMM-based expressive speech...", IEICE Trans. E90-D(9), '07.
- Yoshimura;'01 - "Mixed excitation for HMM-based speech synthesis," Eurospeech, '01.
- Kawahara;'97 - "Restructuring speech representations using a ...", Speech Communication, 27(3), '97.
- Zen;'07 - "Details of the Nitech HMM-based speech synthesis system...", IEICE Trans. E90-D(1), '07.
- Abdl-Hamid;'06 - "Improving Arabic HMM-based speech synthesis quality," Interspeech, '06.

References (3/4)

- Hemptinne;'06 - "Integration of the harmonic plus noise model into the...," Master thesis, IDIAP, '06.
- Banos;'08 - "Flexible harmonic/stochastic modeling...," V. Jornadas en Tecnologias del Habla, '08.
- Cabral;'07 - "Towards an improved modeling of the glottal source in...," ISCA SSW6, '07.
- Maia;'07 - "An excitation model for HMM-based speech synthesis based on ...," ISCA SSW6, '07.
- Ratio;'08 - "HMM-based Finnish text-to-speech system utilizing glottal inverse filtering," Interspeech, '08.
- Drugman;'09 - "Using a pitch-synchronous residual codebook for hybrid HMM/frame...," ICASSP, '09.
- Dines;'01 - "Trainable speech synthesis with trended hidden Markov models," ICASSP, '01.
- Sun;'09 - "Polynomial segment model based statistical parametric speech synthesis...," ICASSP, '09.
- Bulyko;'02 - "Robust splicing costs and efficient search with BMM models for...," ICASSP, '02.
- Shannon;'09 - "Autoregressive HMMs for speech synthesis," Interspeech, '09.
- Zen;'06 - "Reformulating the HMM as a trajectory model...," Computer Speech & Language, 21(1), '06.
- Wu;'06 - "Minimum generation error training for HMM-based speech synthesis," ICASSP, '06.
- Hashimoto;'09 - "A Bayesian approach to HMM-based speech synthesis," ICASSP, '09.
- Wu;'08 - "Minimum generation error training with log spectral distortion for...," Interspeech, '08.
- Toda;'08 - "Statistical approach to vocal tract transfer function estimation based on...," ICASSP, '08.
- Oura;'08 - "Simultaneous acoustic, prosodic, and phrasing model training for TTS...," ISCSLP, '08.
- Ferguson;'80 - "Variable duration models...," Symposium on the application of HMM to text speech, '80.
- Levinson;'86 - "Continuously variable duration hidden...," Computer Speech & Language, 1(1), '86.
- Beal;'03 - "Variational algorithms for approximate Bayesian inference," PhD thesis, Univ. of London, '03.
- Masuko;'03 - "A study on conditional parameter generation from HMM...," Autumn meeting of ASJ, '03.
- Yu;'07 - "A novel HMM-based TTS system using both continuous HMMs and discrete...," ICASSP, '07.

References (4/4)

- Qian;'08 - "Generating natural F0 trajectory with additive trees," Interspeech, '08.
- Latorre;'08 - "Multilevel parametric-base F0 model for speech synthesis," Interspeech, '08.
- Tiomkin;'08 - "Statistical text-to-speech synthesis with improved dynamics," Interspeech, '08.
- Toda;'07 - "A speech parameter generation algorithm considering global...," IEICE Trans. E90-D(5), '07.
- Wu;'08 - "Minimum generation error criterion considering global/local variance...," ICASSP, '08.
- Toda;'09 - "Trajectory training considering global variance for HMM-based speech...," ICASSP, '09.
- Saino;'06 - "An HMM-based singing voice synthesis system," Interspeech, '06.
- Tsuzuki;'04 - "Constructing emotional speech synthesizers with limited speech...," Interspeech, '04.
- Sako;'00 - "HMM-based text-to-audio-visual speech synthesis," ICSLP, '00.
- Tamura;'98 - "Text-to-audio-visual speech synthesis based on parameter generation...," ICASSP, '98.
- Haoka;'02 - "HMM-based synthesis of hand-gesture animation," IEICE Technical report, 102(517), '02.
- Niwase;'05 - "Human walking motion synthesis with desired pace and...," IEICE Trans. E88-D(11), '05.
- Hofer;'07 - "Speech driven head motion synthesis based on a trajectory model," SIGGRAPH, '07.
- Ma;'07 - "A MSD-HMM approach to pen trajectory modeling for online handwriting...," ICDAR, '07.
- Morioka;'04 - "Miniaturization of HMM-based speech synthesis," Autumn meeting of ASJ, '04.
- Kim;'06 - "HMM-Based Korean speech synthesis system for...," IEEE Trans. Consumer Elec., 52(4), '06.
- Klatt;'82 - "The Klatt-Talk text-to-speech system," ICASSP, '82.