

Statistical approach to speech synthesis ---past, present, and future

Keiichi Tokuda

Nagoya Institute of Technology

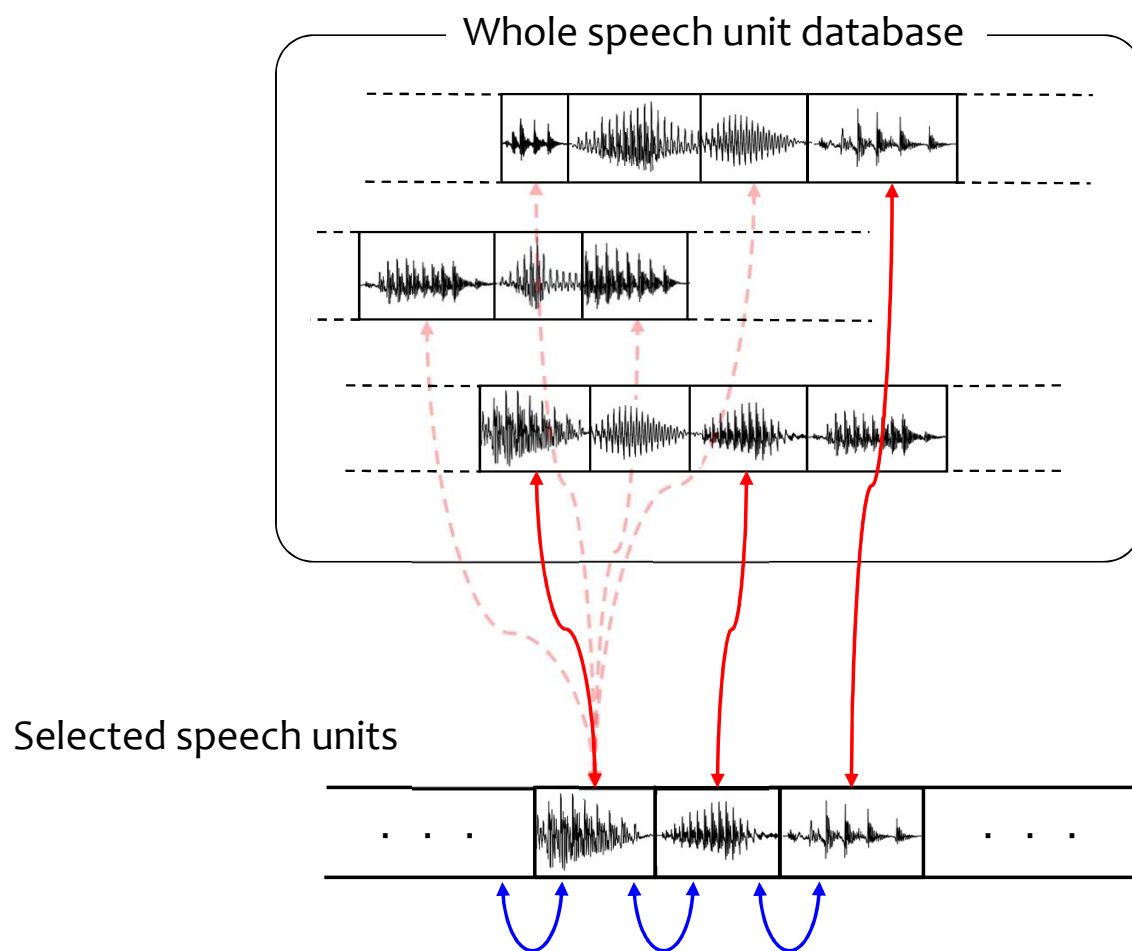
INTERSPEECH 2019

Speech synthesis approaches

- Rule-based, formant synthesis (~'90s)
Phonetic units are built by hand-crafted rules
- Corpus-based, concatenative synthesis ('90s~)
Concatenate speech units (in acoustic feature or waveform) from a database
 - Single inventory: diphone synthesis
 - Multiple inventory: **unit selection synthesis** →
- Corpus-based, statistical synthesis (late '90s~) ←
Source-filter model + statistical acoustic model
 - **HMM** (hidden Markov model) (1995~) >>
 - **DNN** (deep neural networks) (2013~)
 - **WaveNet** (2016~)

We were
working on this

Unit-selection synthesis



When we are lucky:



When we are unlucky:

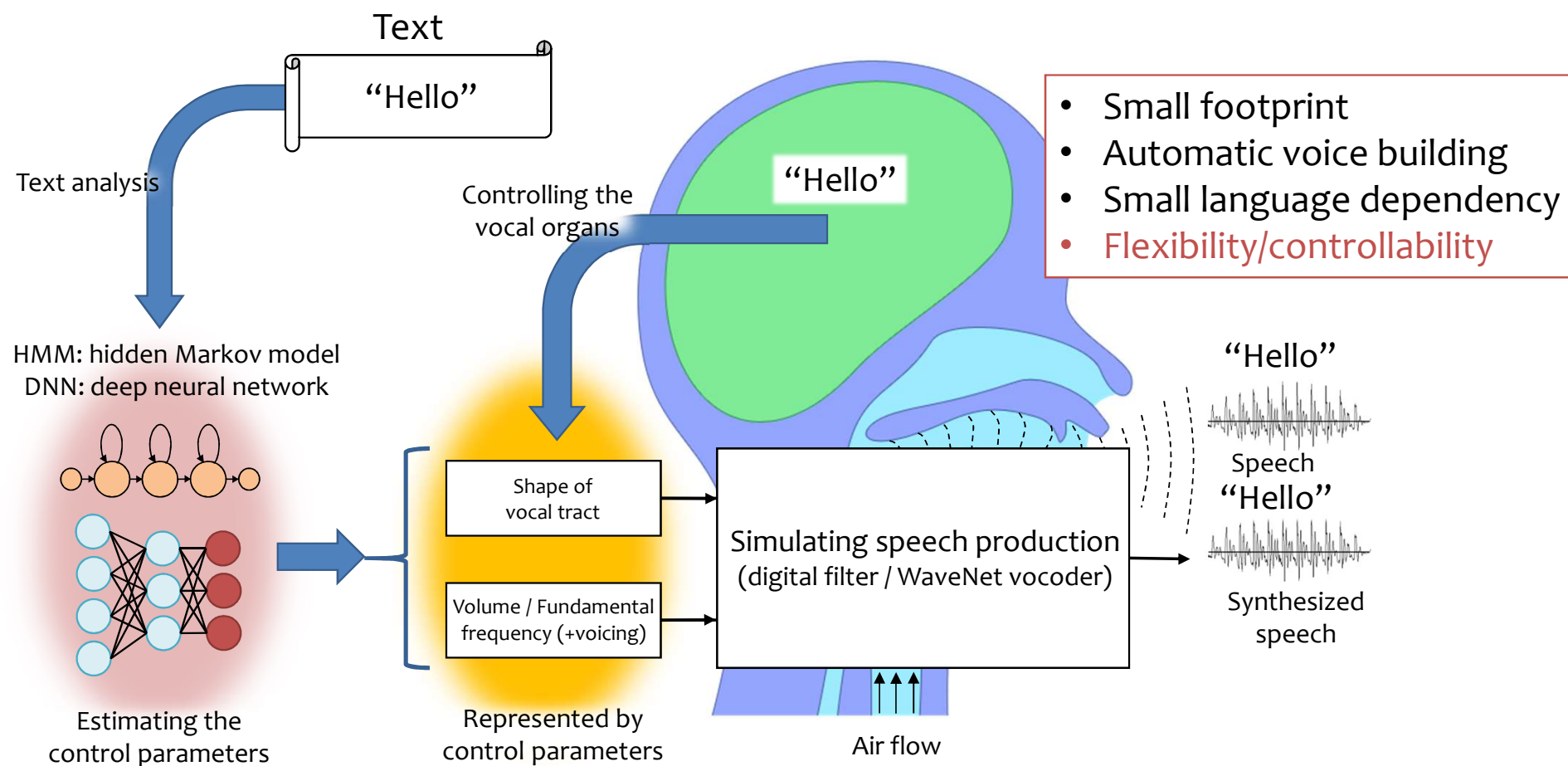


— Target cost

— Concatenation cost

Minimize the total cost in runtime
using dynamic programming

Statistical approach to speech synthesis



Pros and cons

Solved by neural vocoding, e.g., WaveNet

Unit selection	Statistical parametric	
Waveform concatenation → Natural sounding	Vocoded → buzzy or muffled	☹️
Discontinuity, hit or miss	Smooth	
Work better for larger databases	Can work for small databases	
Large footprint	Small footprint	
Fixed voice → fixed style, fixed emotional expression, etc.	Flexible → speaker adaptation, speaking style interpolation, etc.	😊

No reason to hesitate to move onto the statistical approach

Outline

- Statistical formulation of speech synthesis
- HMM-based speech synthesis
- Deep neural networks
- Evaluation / data & software tools
- Other related topics

Outline

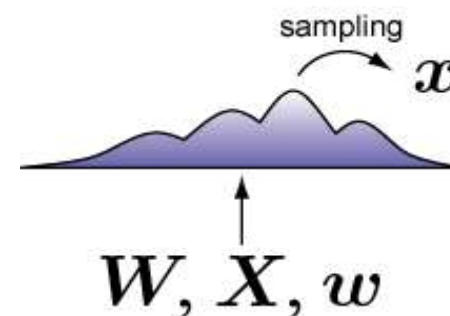
- Statistical formulation of speech synthesis
- HMM-based speech synthesis
- Deep neural networks
- Evaluation / data & software tools
- Other related topics

The basic problem of speech synthesis

We have a speech database, i.e.,
a set of pairs of texts and corresponding speech waveforms.
Given a text to be synthesized,
what is the speech waveform corresponding to the text?

- W : texts
 - X : speech waveforms
- } database
- } Given
- w : text to be synthesized ($w \notin W$)
 - x : speech waveform ← ?

$$x \sim p(x|w, X, W)$$



Statistical formulation of speech synthesis (1/4)

- Estimating predictive distribution is hard.
→ Introduce generative representation (λ : model parameters)

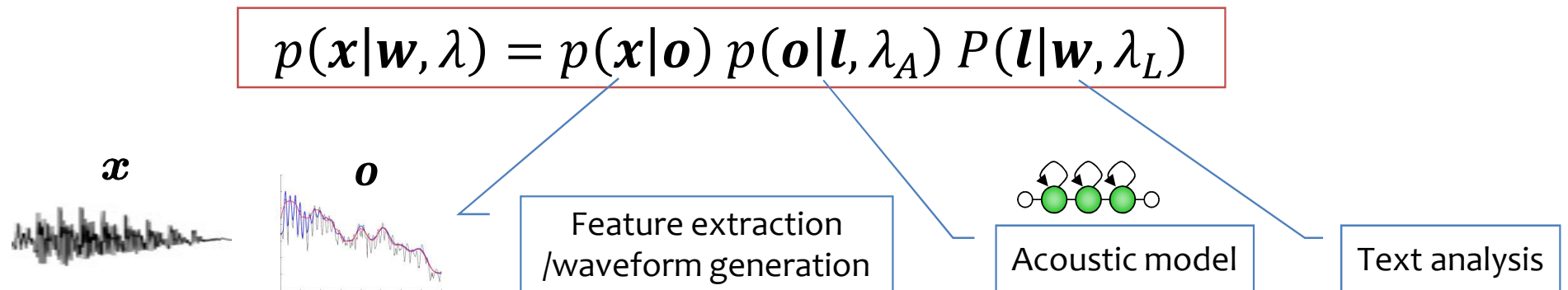
$$p(\mathbf{x}|\mathbf{w}, \mathbf{X}, \mathbf{W}) = \int p(\mathbf{x}|\mathbf{w}, \lambda) p(\lambda|\mathbf{X}, \mathbf{W}) d\lambda$$

- It is difficult to perform integral over auxiliary variables
→ Approximate integral by maximizing $p(\lambda|\mathbf{X}, \mathbf{W})$

$$\begin{aligned} \hat{\lambda} &= \arg \max_{\lambda} p(\hat{\lambda}|\mathbf{W}, \mathbf{X}) \leftarrow \text{training} \\ \mathbf{x} &\sim p(\mathbf{x}|\mathbf{w}, \hat{\lambda}) \leftarrow \text{generation} \end{aligned}$$

Statistical formulation of speech synthesis (2/4)

- Usually the generative model is decomposed into sub-modules, e.g.,



\mathbf{o} : parametric representation of speech waveform \mathbf{x}

\mathbf{l} : linguistic feature for \mathbf{w}

$\lambda = \{\lambda_A, \lambda_L\}$: generative model parameter

λ_A : acoustic model parameter

λ_L : text analysis module parameter

Linguistic feature

Phoneme (or distinctive feature)

- {preceding, current, succeeding} phonemes

Syllable

- # of phonemes in {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, etc.

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
- Syntactic/dependency information, etc.

(→)

Phrase

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

⋮

+Frame-level Duration and positional information

+Speaking styles, emotional expressions, etc.
when we have such tags/labels

Statistical formulation of speech synthesis (3/4)

- Decompose the generative model into sub-modules:

$$\begin{aligned}\hat{\lambda} &= \arg \max_{\lambda} p(\hat{\lambda} | \mathbf{W}, \mathbf{X}) \leftarrow \text{training} \\ \mathbf{x} &\sim p(\mathbf{x} | \mathbf{w}, \hat{\lambda}) \leftarrow \text{generation}\end{aligned}$$



$$\begin{aligned}\{\hat{\lambda}_A, \hat{\lambda}_L\} &= \arg \max_{\lambda_A, \lambda_L} \int \sum_L p(\mathbf{X} | \mathbf{o}) p(\mathbf{o} | \mathbf{L}, \lambda_A) P(\mathbf{L} | \mathbf{W}, \lambda_L) d\mathbf{o} p(\lambda_A) p(\lambda_L) \\ \mathbf{x} &\sim \int \sum_l p(\mathbf{x} | \mathbf{o}) p(\mathbf{o} | \mathbf{l}, \hat{\lambda}_A) P(\mathbf{l} | \mathbf{w}, \hat{\lambda}_L) d\mathbf{o} \leftarrow \text{generation}\end{aligned}$$

\uparrow training

Statistical formulation of speech synthesis (4/4)

- It is difficult to perform integral and sum
→ Approximated by step-by-step maximization

←

$\hat{\lambda}_L$: pre-trained text analysis module parameter

$\hat{\mathbf{O}} = \arg \max_{\mathbf{O}} p(\mathbf{X}|\mathbf{O}) \leftarrow$ speech feature parameter extraction

$\hat{\mathbf{L}} = \arg \max_{\mathbf{L}} P(\mathbf{L}|\mathbf{W}, \hat{\lambda}_L)$ or $p(\hat{\mathbf{O}}|\mathbf{L}, \hat{\lambda}_A)$ or $p(\hat{\mathbf{O}}|\mathbf{L}, \hat{\lambda}_A)P(\mathbf{L}|\mathbf{W}, \hat{\lambda}_L) \leftarrow$ labeling

$\hat{\lambda}_A = \arg \max_{\lambda_A} p(\hat{\mathbf{O}}|\hat{\mathbf{L}}, \lambda_A)p(\lambda_A) \leftarrow$ acoustic model training

↑ Training

$\hat{\mathbf{l}} = \arg \max_{\mathbf{l}} P(\mathbf{l}|\mathbf{w}, \hat{\lambda}_L) \leftarrow$ text analysis

↓ Synthesis

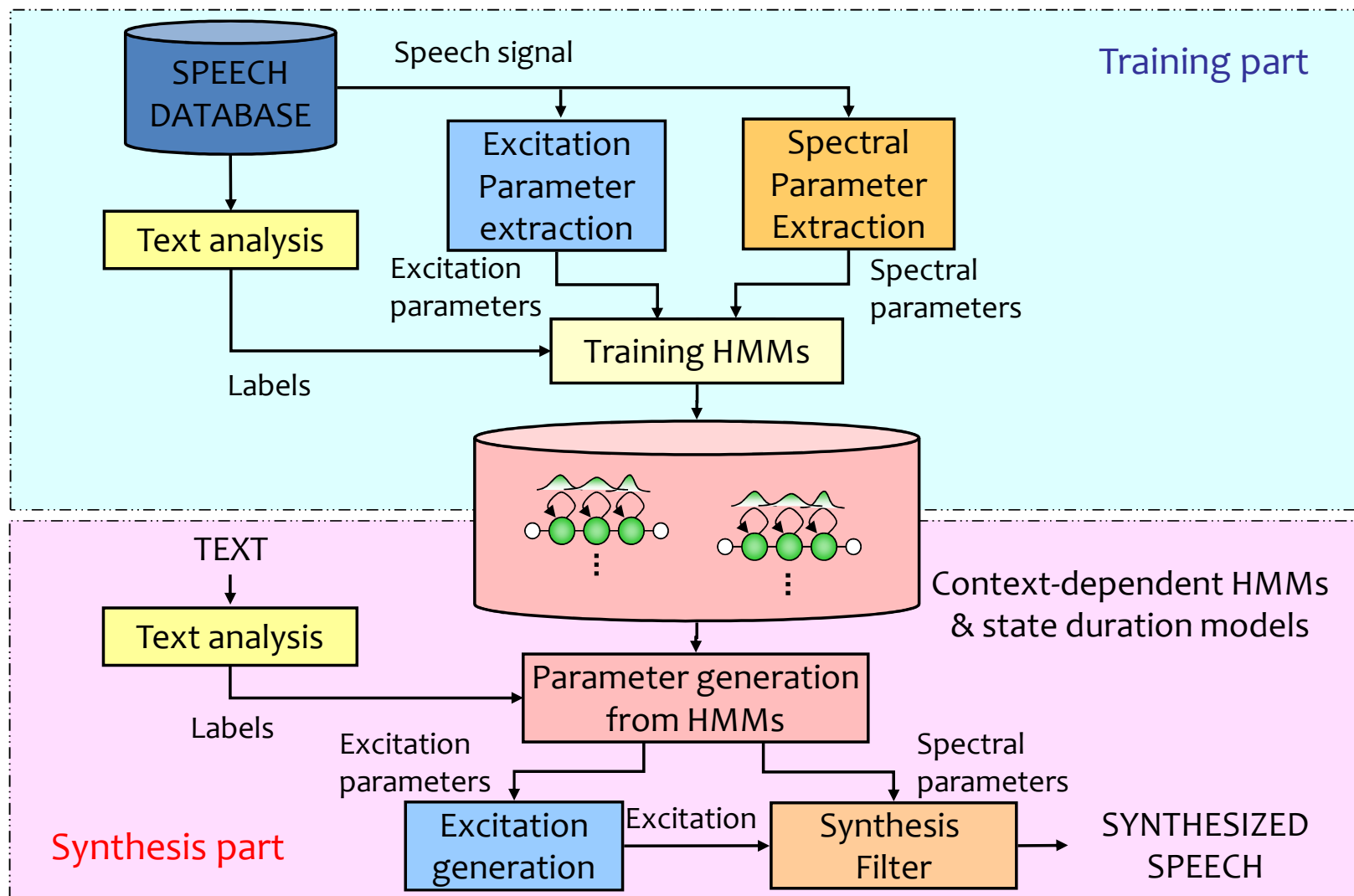
$\hat{\mathbf{o}} = \arg \max_{\mathbf{o}} p(\mathbf{o}|\hat{\mathbf{l}}, \hat{\lambda}_A) \leftarrow$ speech parameter generation

$\mathbf{x} \sim p(\mathbf{x}|\hat{\mathbf{o}}) \leftarrow$ waveform generation

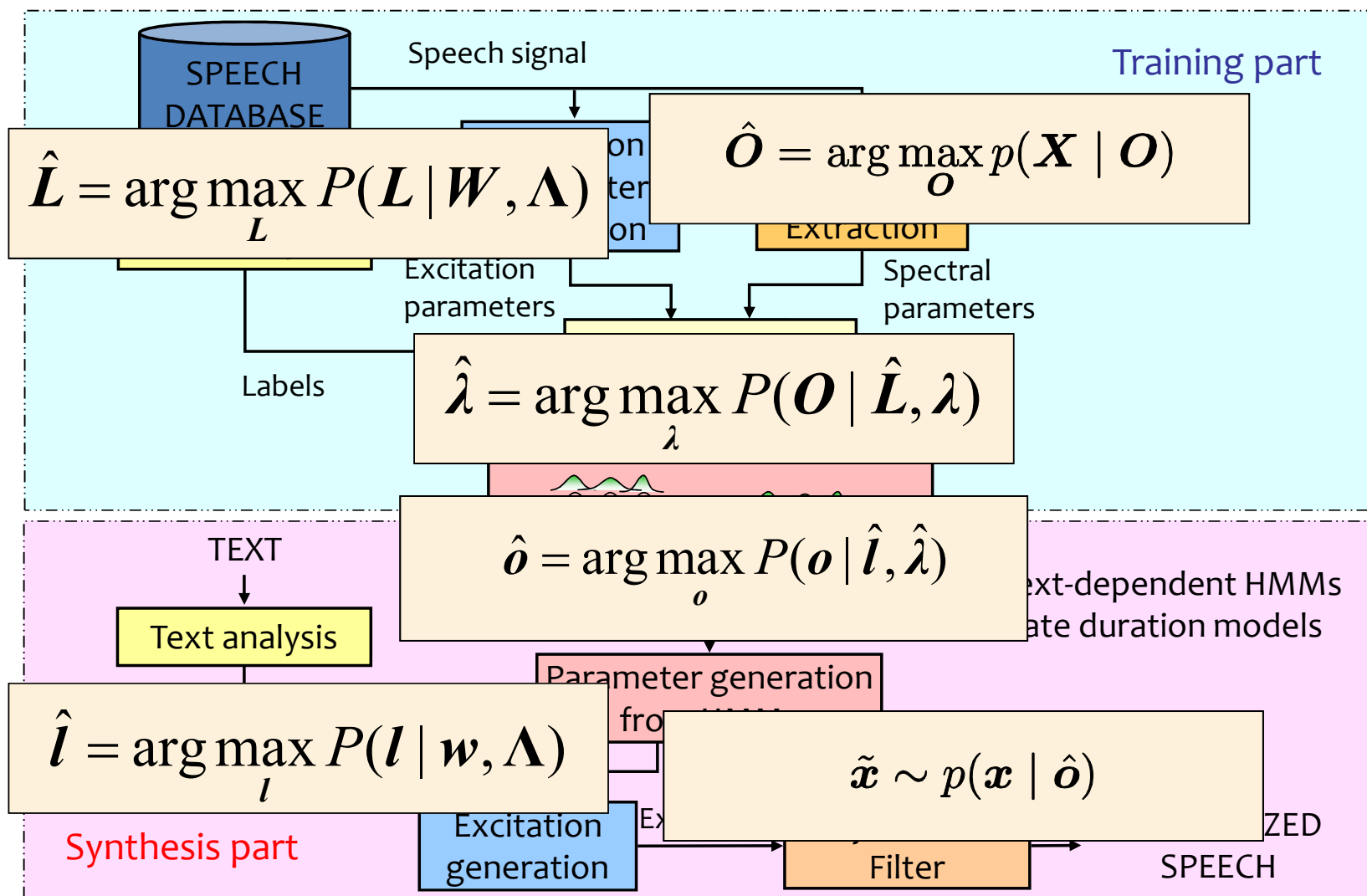
Outline

- Statistical formulation of speech synthesis
- **HMM-based speech synthesis**
- Deep neural networks
- Evaluation / data & software tools
- Other related topics

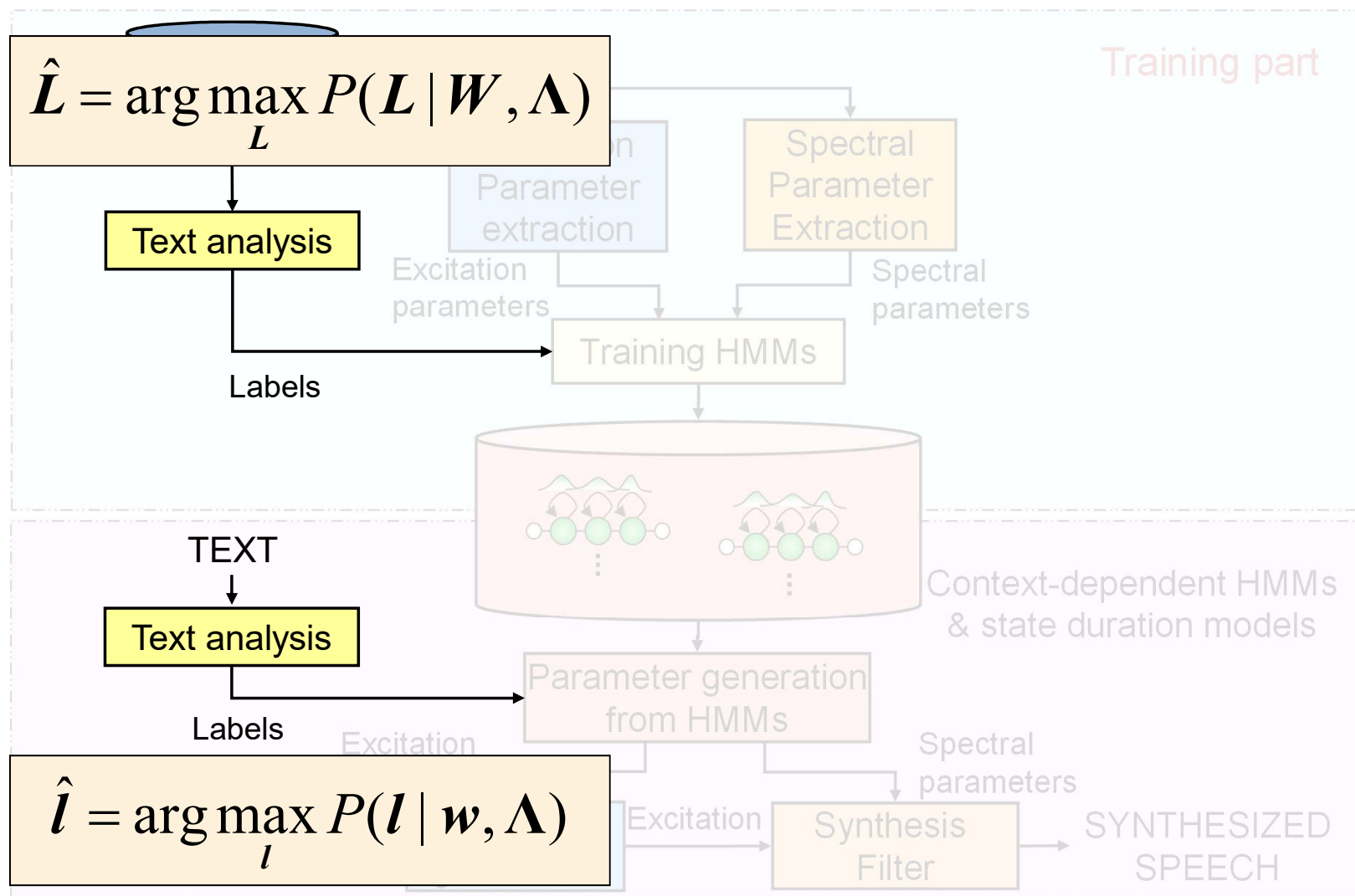
HMM-based speech synthesis system



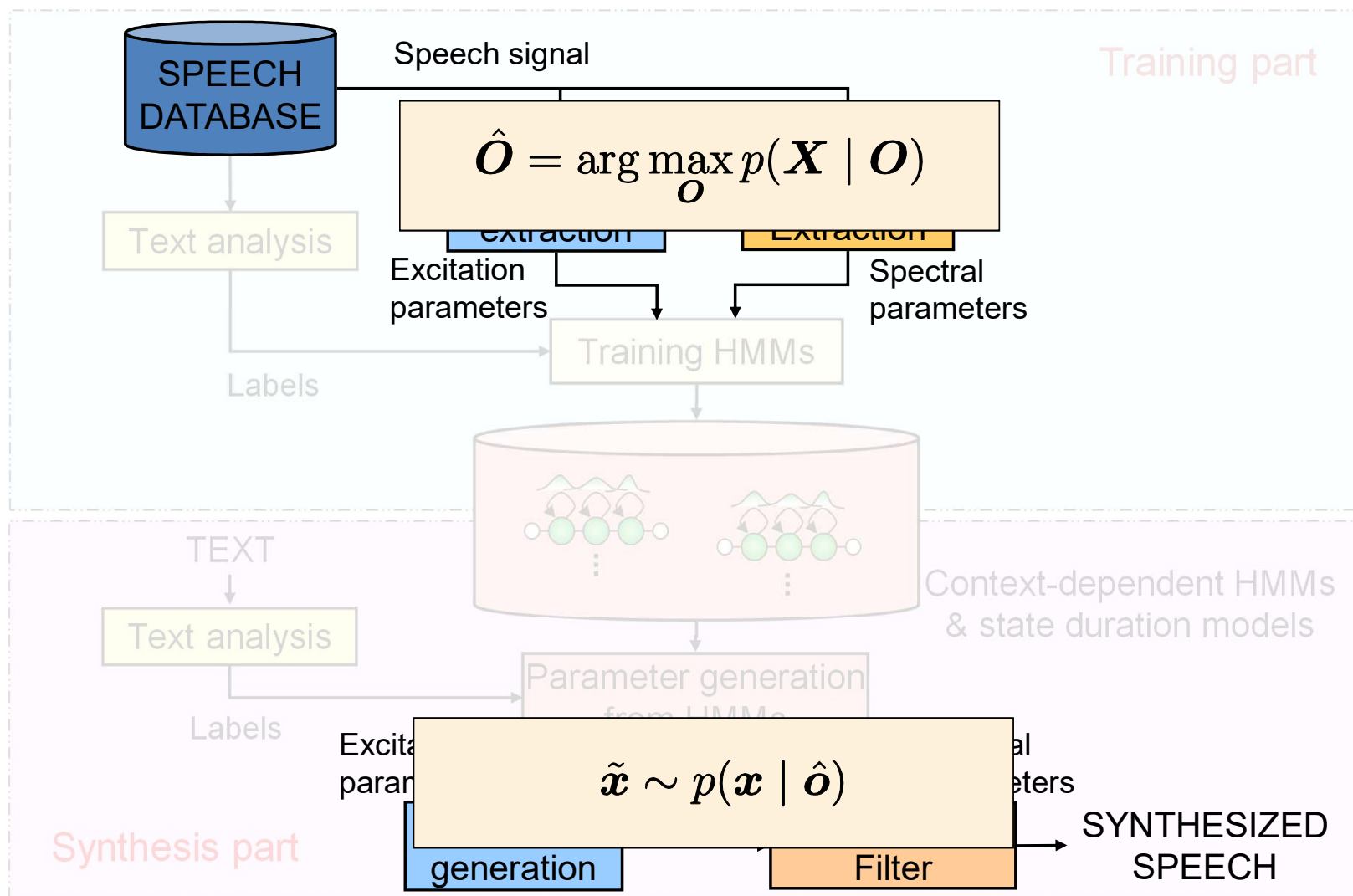
HMM-based speech synthesis system



HMM-based speech synthesis system



HMM-based speech synthesis system



Mel-cepstrum-based spectral analysis

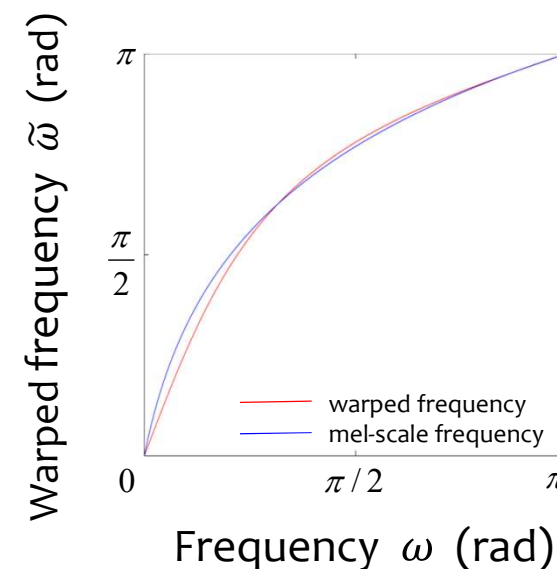
$$H(e^{j\omega}) = \exp \sum_{m=0}^M c(m) e^{-j\tilde{\omega}m}, \quad e^{-j\tilde{\omega}} = \frac{e^{-j\omega} - \alpha}{1 - \alpha e^{-j\omega}}$$

$$\mathbf{c} = [c(0), c(1), \dots, c(M)]^T \leftarrow \text{mel-cepstrum}$$

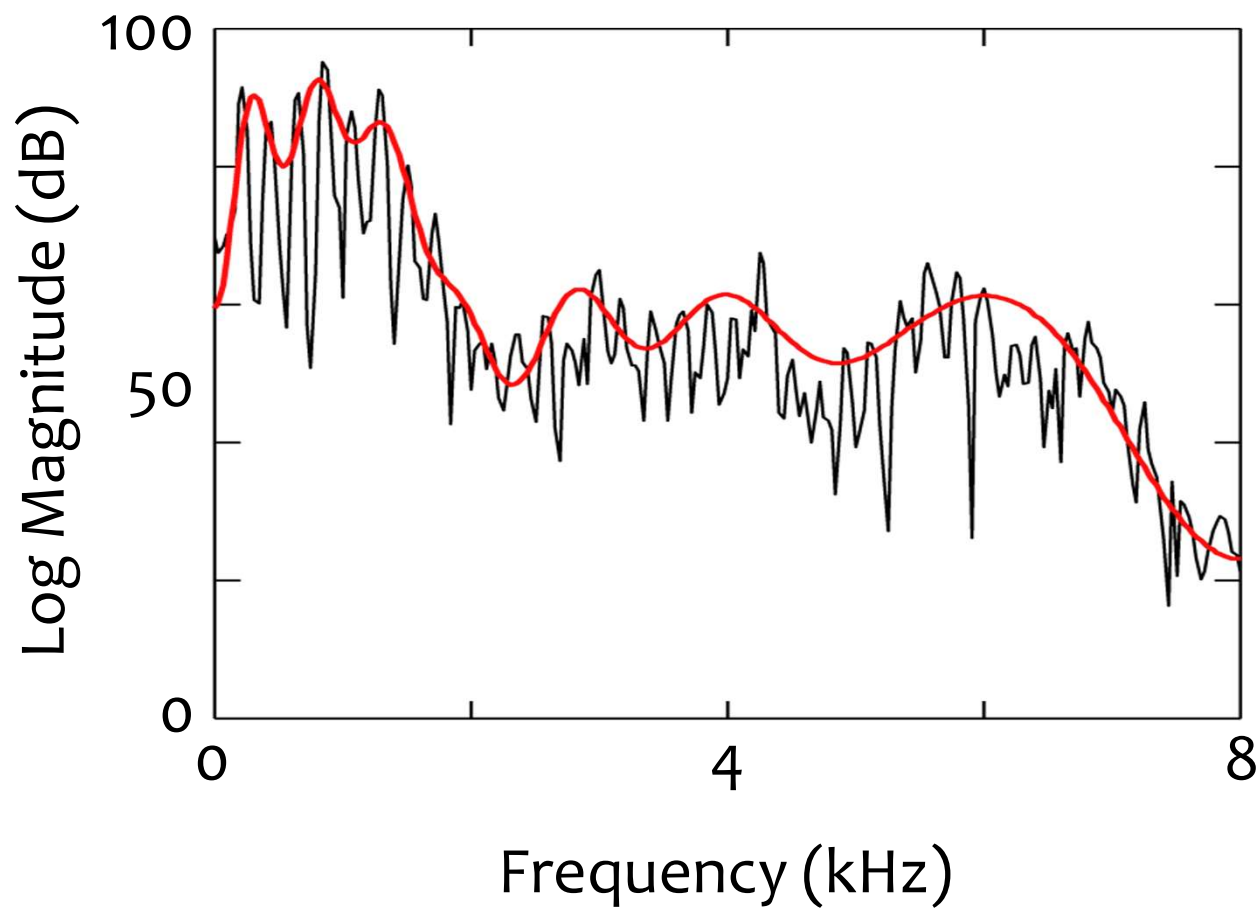
ML estimation of mel-cepstrum:

$$\hat{\mathbf{c}} = \arg \max_{\mathbf{c}} p(\mathbf{x}|\mathbf{c}) \leftarrow$$

when \mathbf{x} is Gaussian process,
 $p(\mathbf{x}|\mathbf{c})$ is convex with respect to \mathbf{c} [Fuka92]

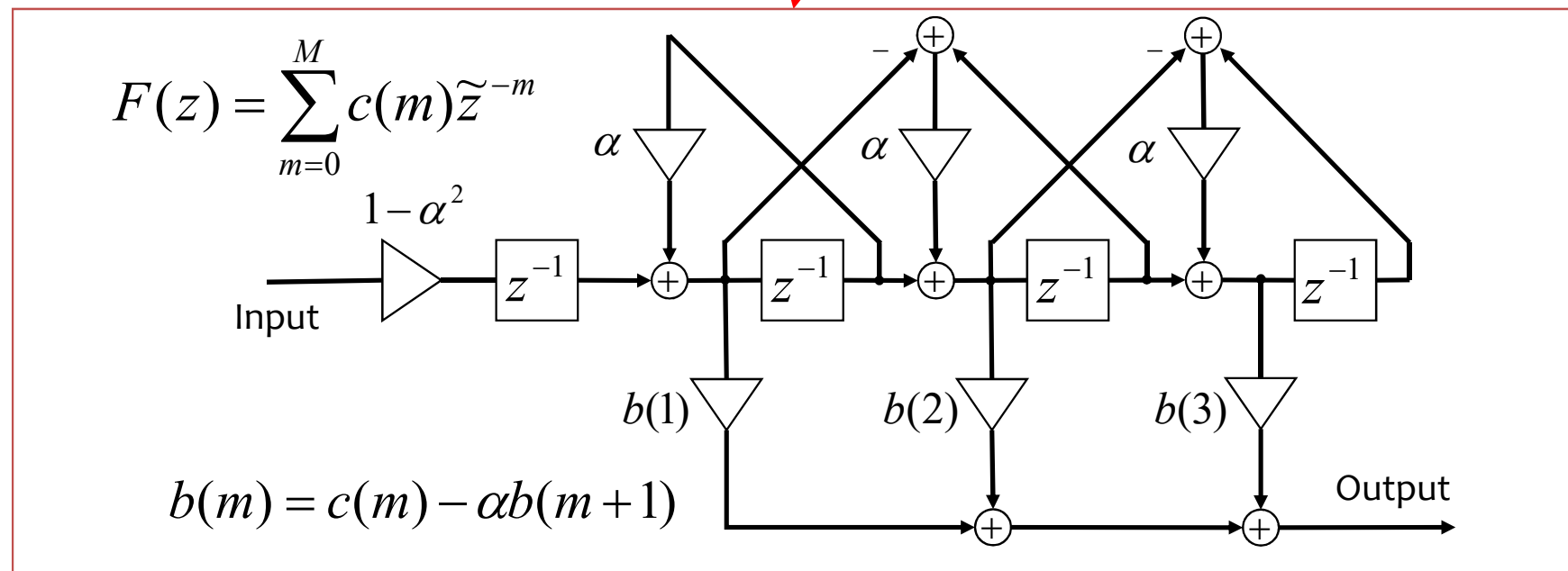


Spectral estimation example



MLSA filter (1/2) [Fukada '92]

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

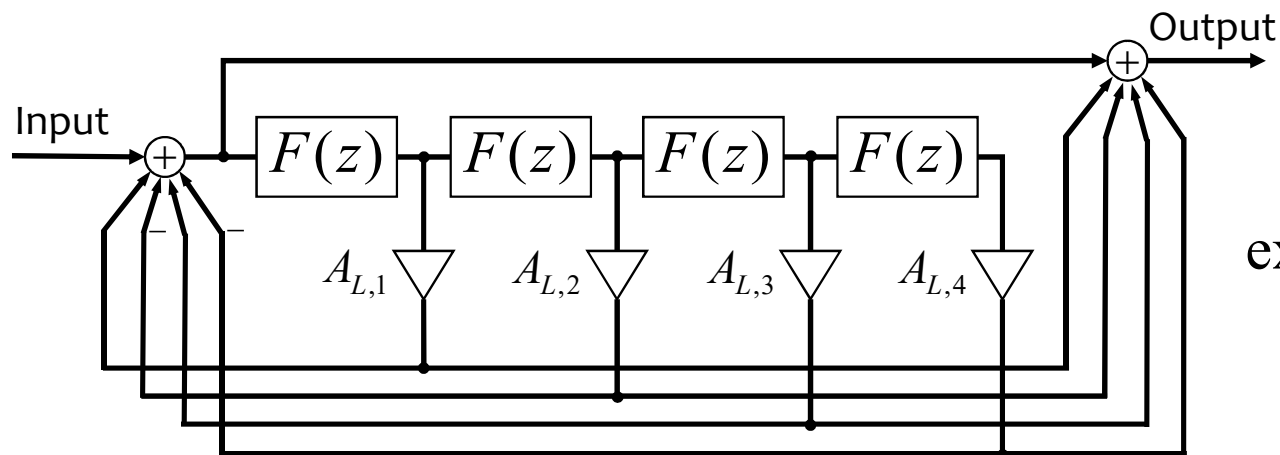


MLSA filter (2/2) [Fukada '92]

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

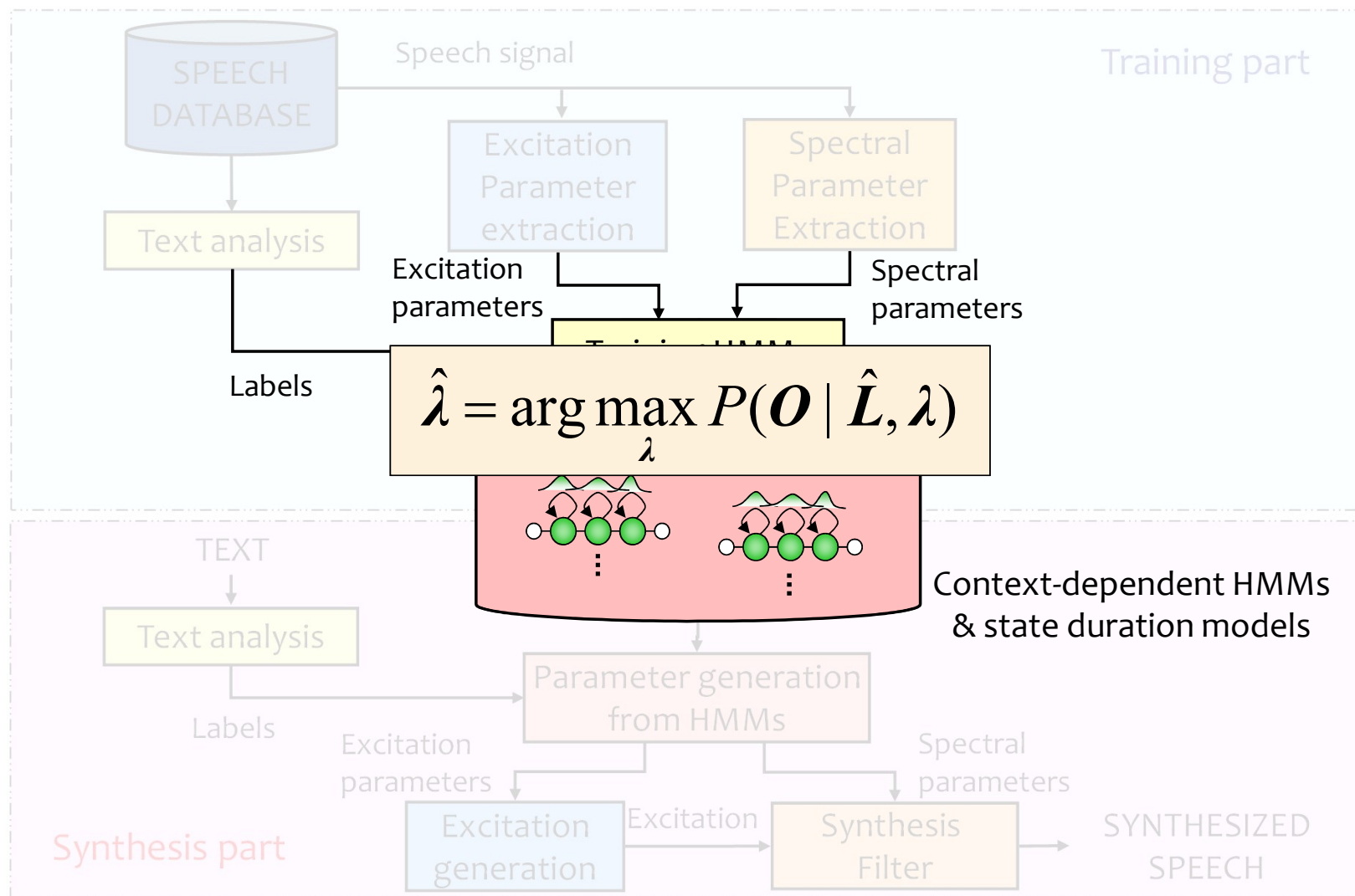
$$\exp x \cong \frac{1 + \sum_{l=1}^L A_{L,l} x^l}{1 + \sum_{l=1}^L A_{L,l} (-x)^l}$$

- Approximation error < 0.24dB
- O(8M) operations/sample
- Stable filter

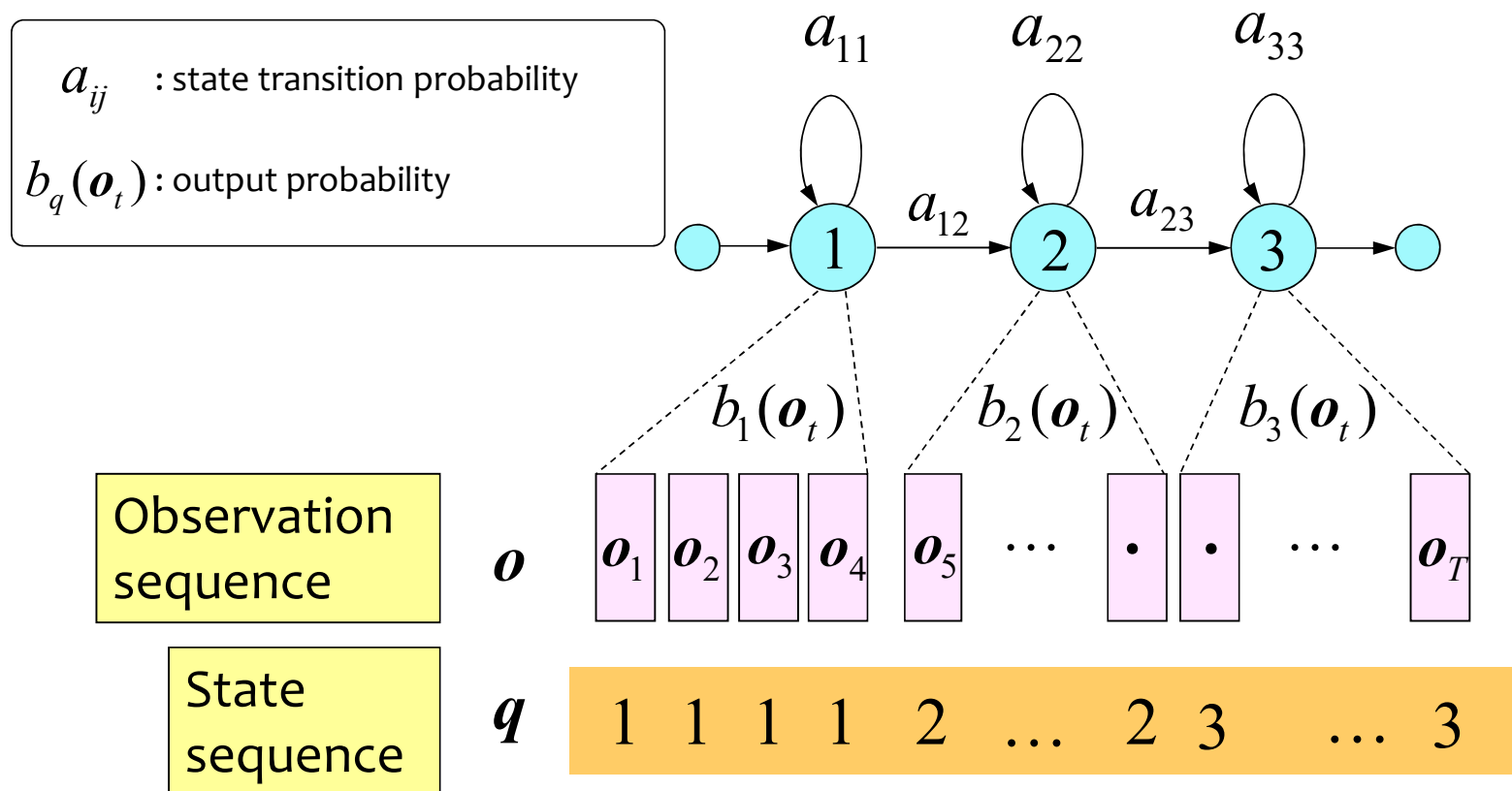


$$\exp F(z) \cong \frac{1 + \sum_{l=1}^L A_{L,l} \{F(z)\}^l}{1 + \sum_{l=1}^L A_{L,l} \{-F(z)\}^l}$$

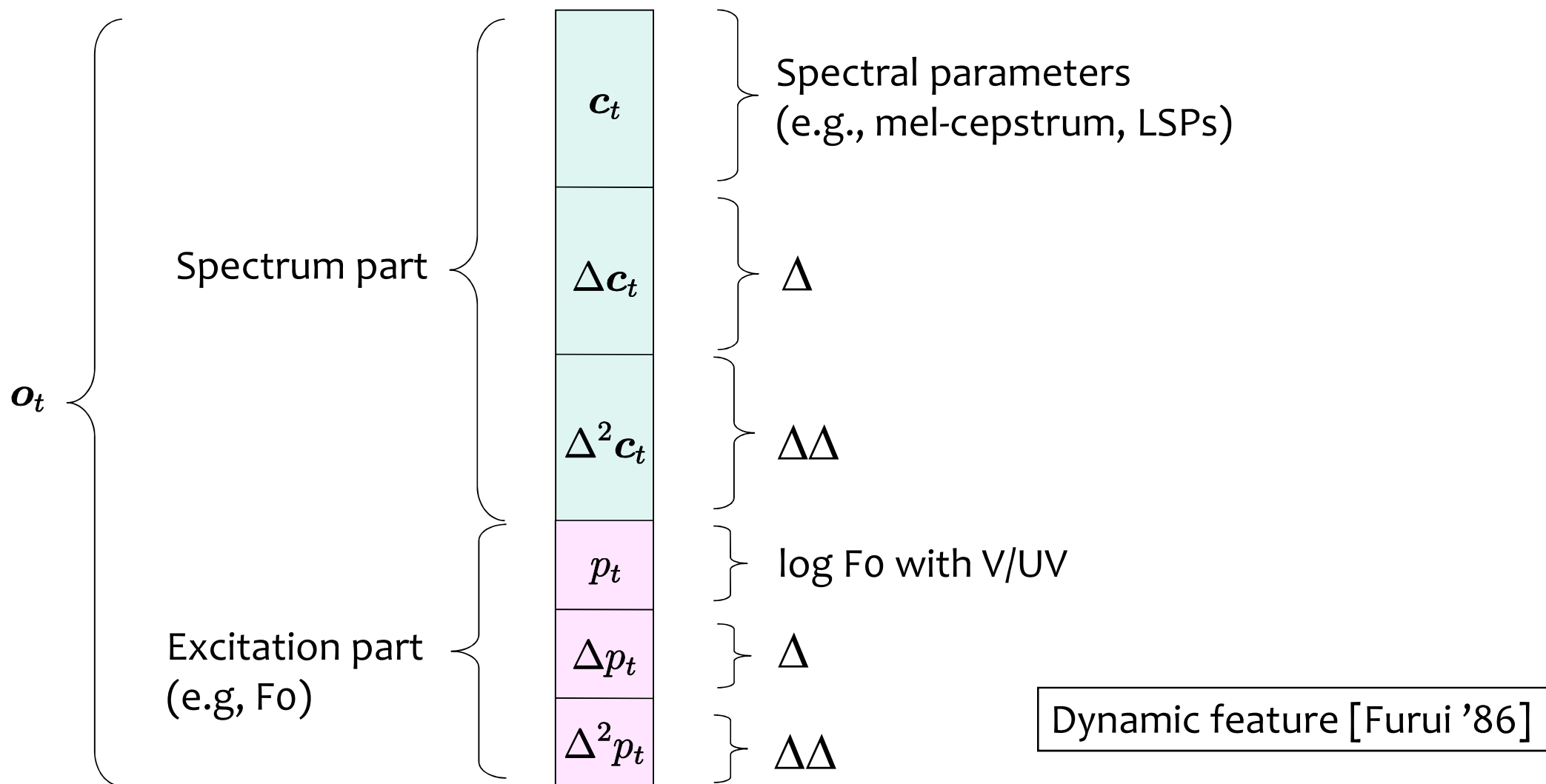
HMM-based speech synthesis system



Hidden Markov model (HMM)

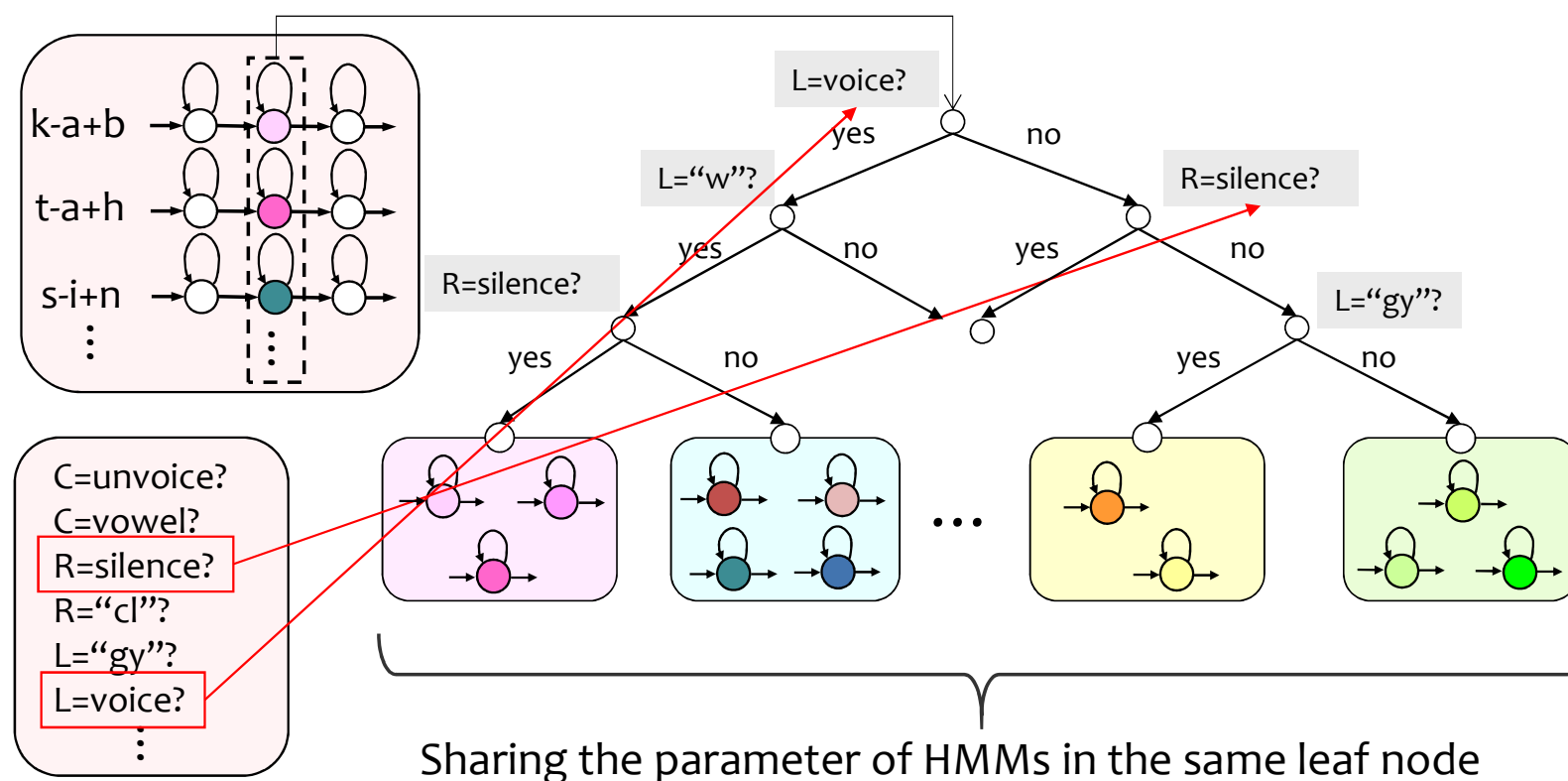


Structure of state output (observation) vector



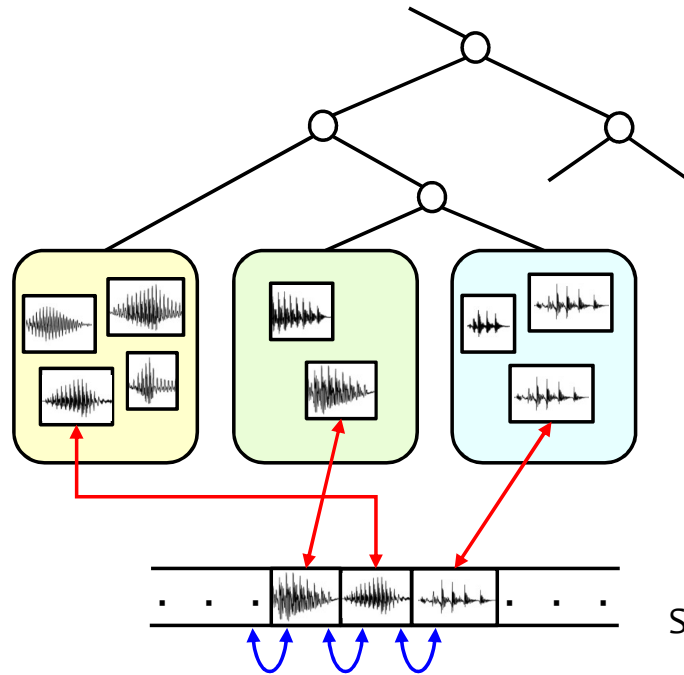
Decision tree-based state clustering [Odell; '95]

$p(\mathbf{o}|\mathbf{l}, \lambda_A) : \text{HMM}$, \mathbf{l} : linguistic feature

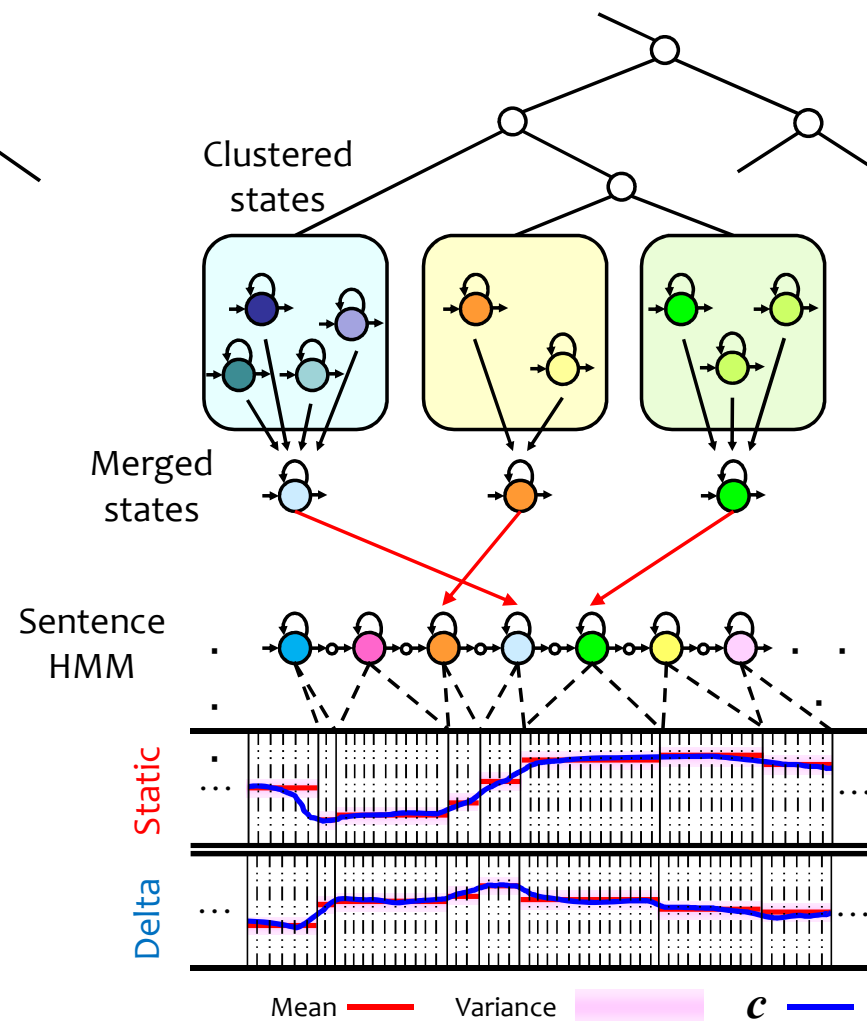


Relation between two approaches

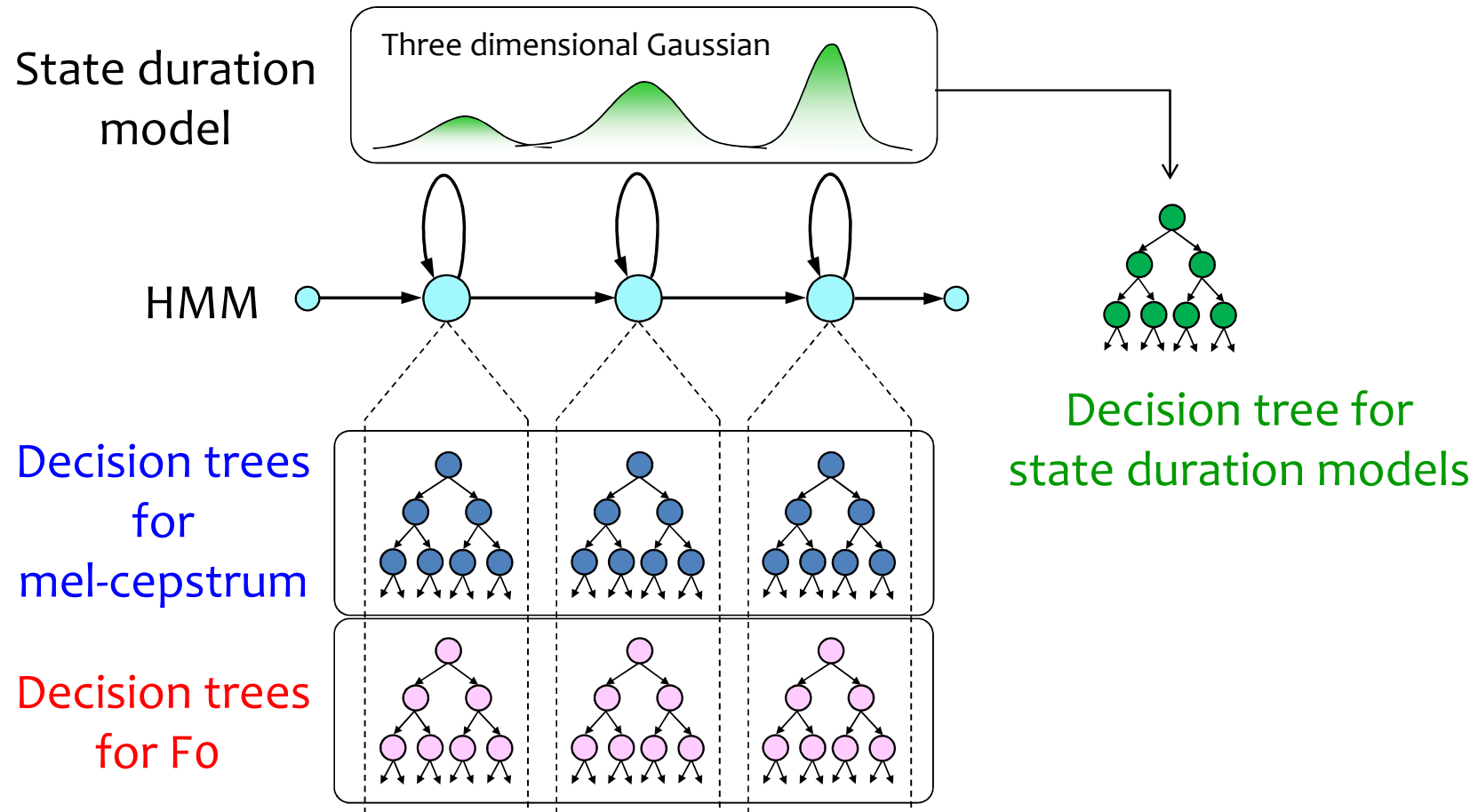
Unit Selection



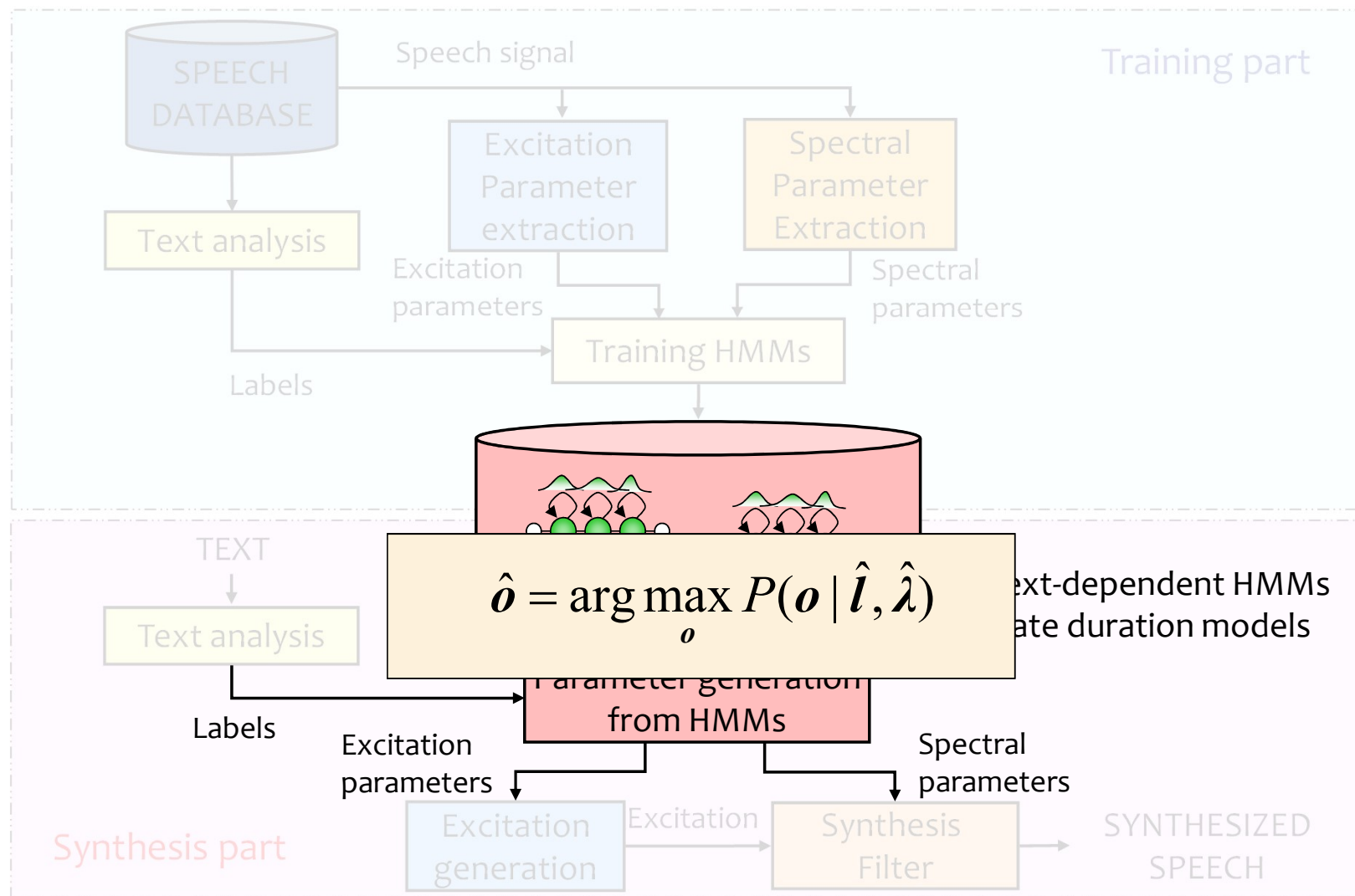
HMM-based



Stream-dependent tree-based clustering



HMM-based speech synthesis system



Speech parameter generation algorithm

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

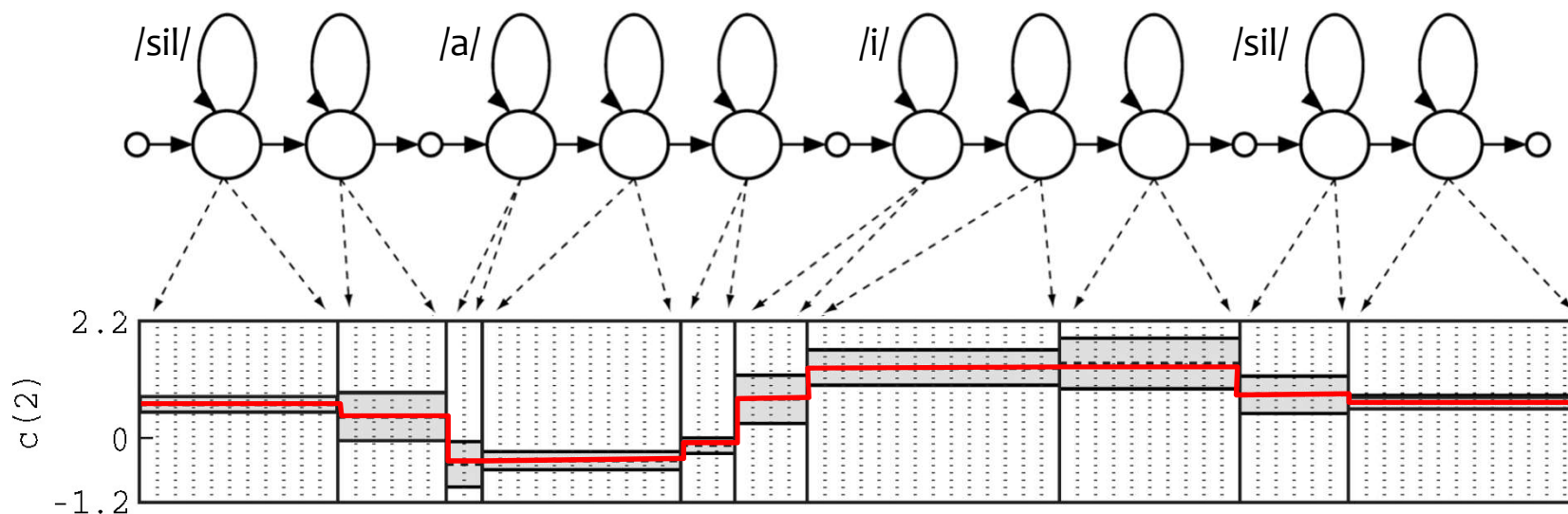
q : state sequence



$$\hat{q} = \arg \max_q P(q | \hat{l}, \hat{\lambda}) \quad \leftarrow \text{Determination of durations}$$

$$\hat{o} = \arg \max_o P(o | \hat{q}, \hat{\lambda}) \quad \leftarrow \text{Determination of speech parameter}$$

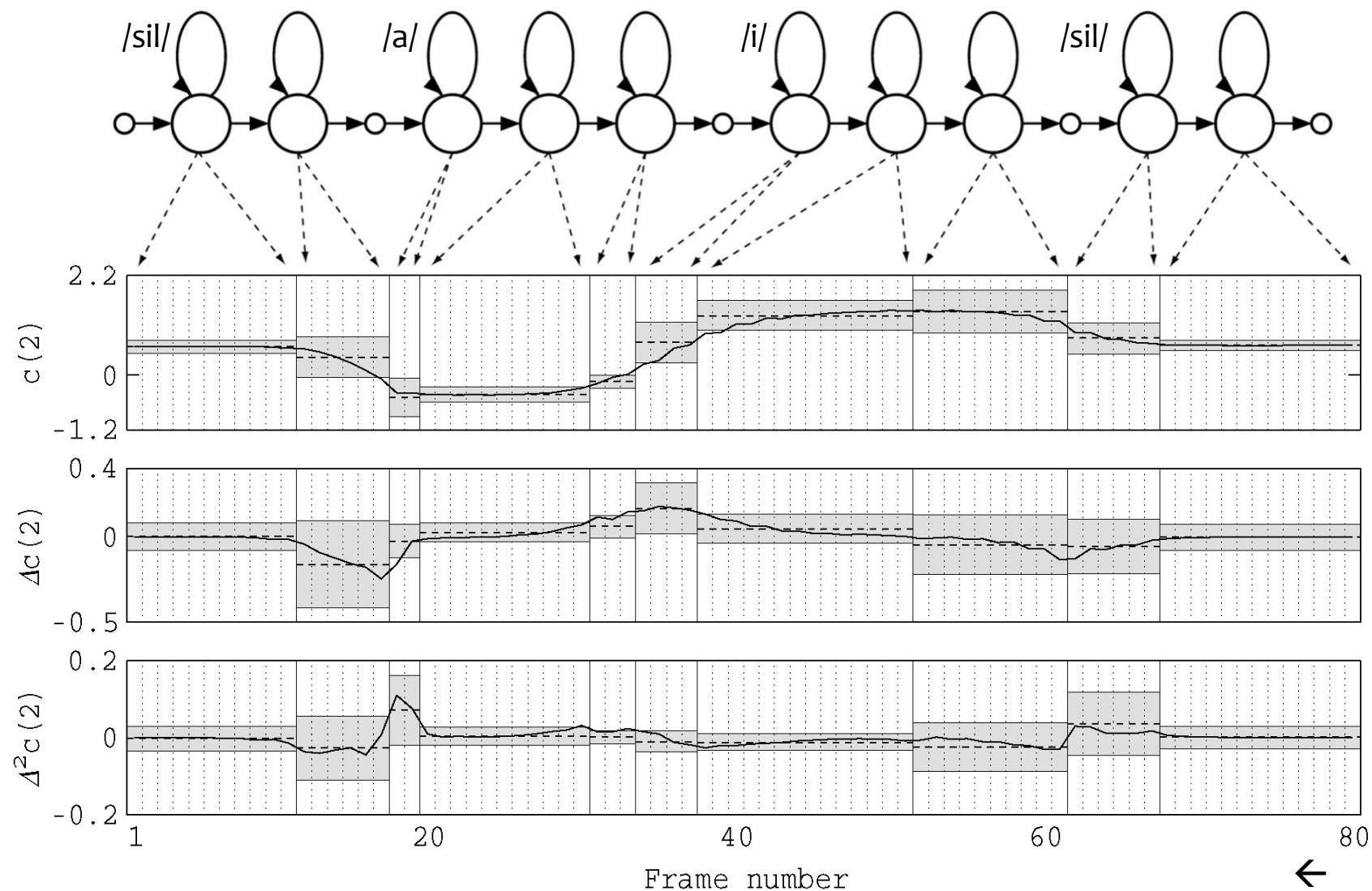
Generated speech parameter trajectory



Frame number



Generated speech parameter trajectory



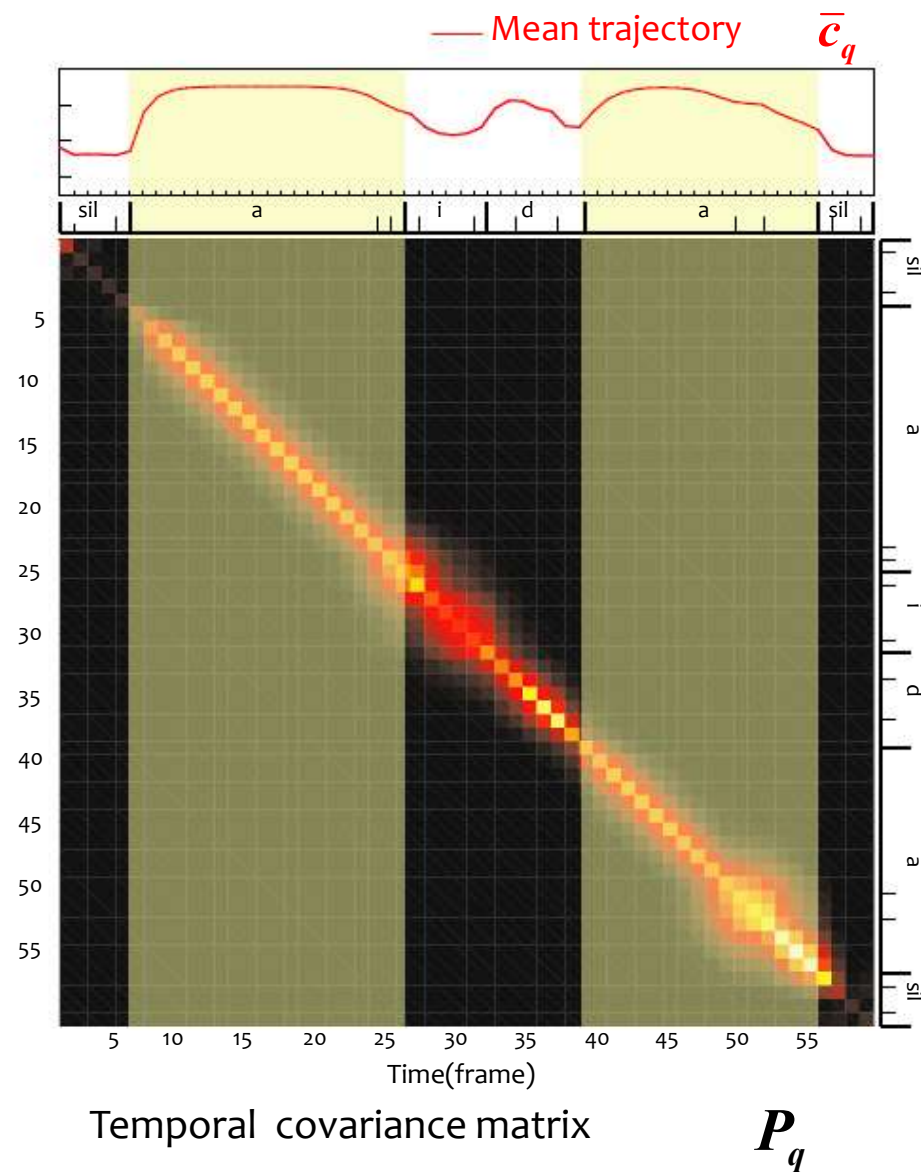
Trajectory HMM

with dynamic feature

w/o dynamic feature

$$\frac{1}{Z_c} P(\mathbf{o}|\mathbf{q}, \hat{\lambda}) = N(\mathbf{c}|\bar{\mathbf{c}}_q, \mathbf{P}_q)$$

$$Z_c = \int P(\mathbf{o}|\mathbf{q}, \hat{\lambda}) d\mathbf{c}$$



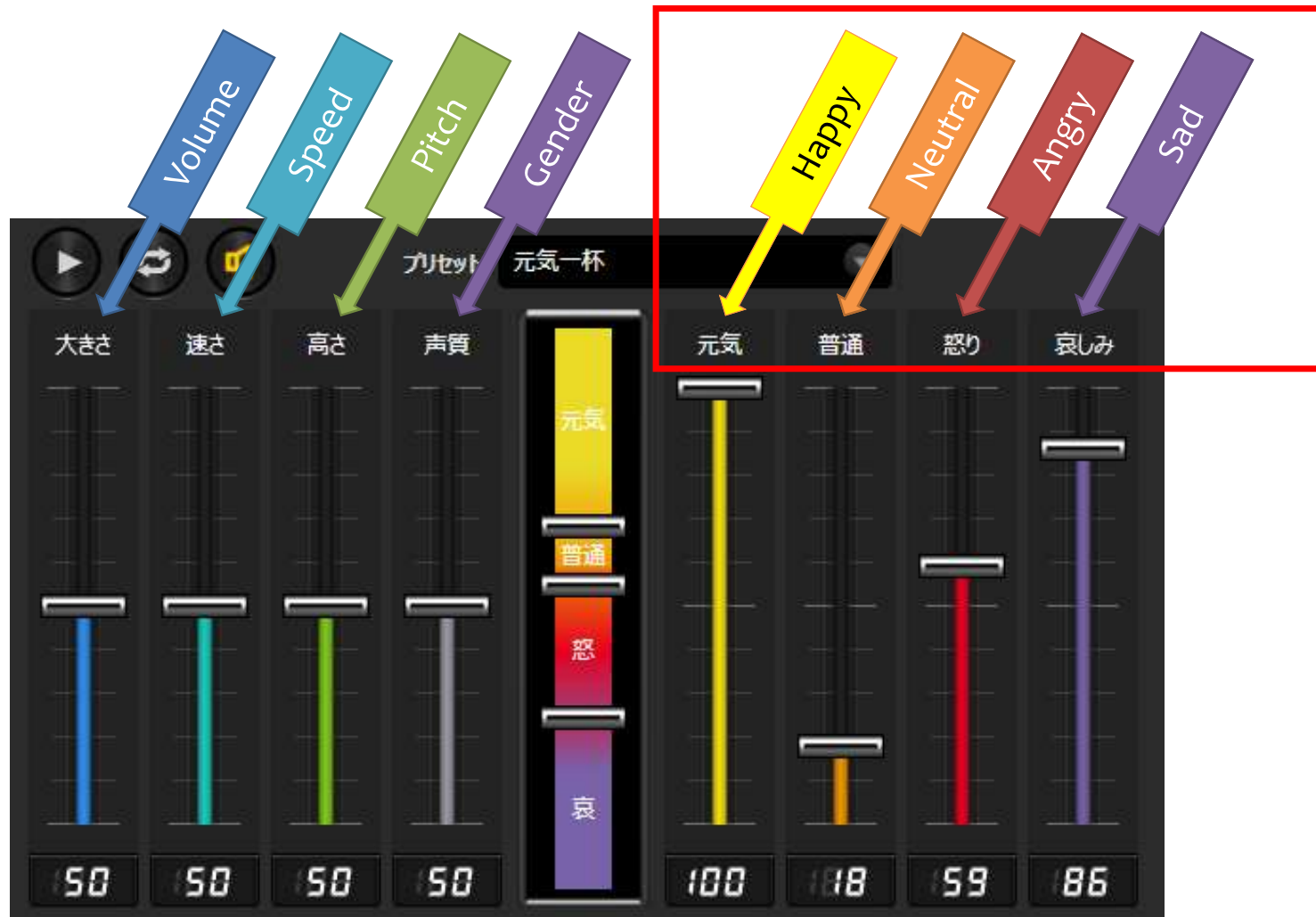
Flexibility to control speech variations

- **Speaker Adaptation** (mimicking voices)
 - [Tamura '98], [Tamura '01], [Yamagishi '03], ...
- **Speaker Interpolation** (mixing voices)
 - [Yoshimura '97], ...
- **Eigenvoice** (producing voices)
 - [Shichiri '02], [Kazumi '10], ...
- **Multiple-regression** (controlling voices)
 - [Nose '07], ...

(→)

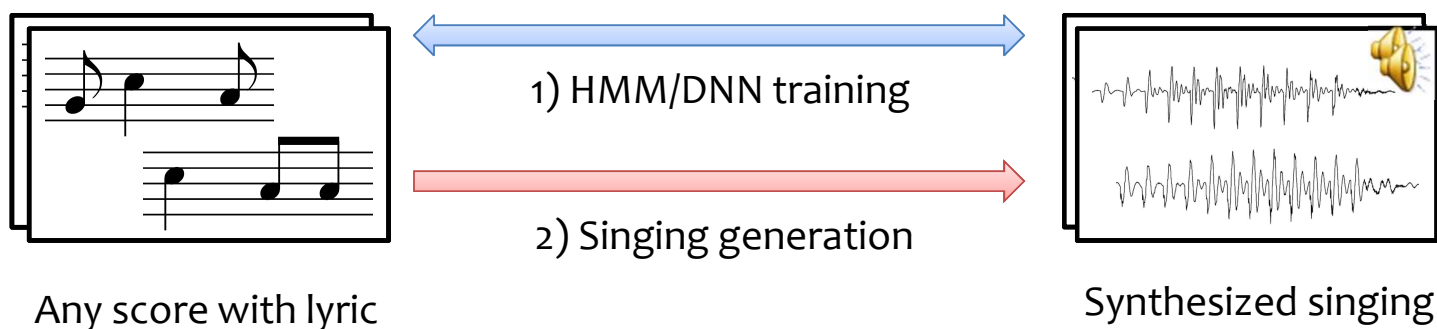
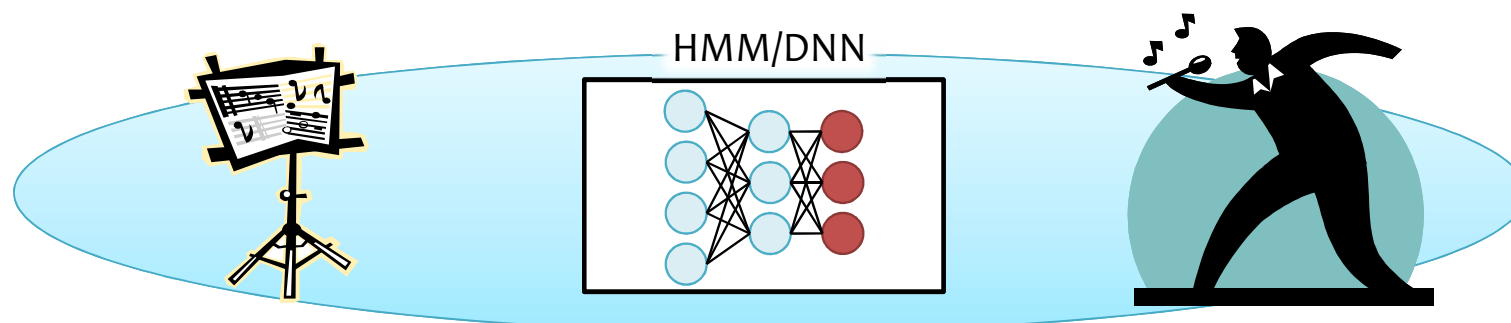
Only from publications by the HTS working group

Mixing emotional expressions



(←)

Singing synthesis



HMM+STRAIGHT

Outline

- Statistical formulation of speech synthesis
- HMM-based speech synthesis
- **Deep neural networks**
- Evaluation / data & software tools
- Other related topics

Hidden Markov model approach

STRAIGHT HMM

$$p(\mathbf{x}|\mathbf{w}, \lambda) = p(\mathbf{x}|\mathbf{c}) p(\mathbf{c}|\mathbf{l}, \lambda_A) P(\mathbf{l}|\mathbf{w}, \lambda_L)$$

- GV-based parameter generation [Toda '05]
- HSMM (hidden semi-Markov model) [Zen '07]
- Trajectory HMM training [Zen '07]
- MGE training [Wu '08]
- Bayesian approach [Hashimoto '09]
- Additive decision tree [Takaki '10]
- Trainable excitation model [Maia '07], etc.

Text analysis

Only from publications by the HTS working group

Recombining submodules

$$p(\mathbf{x}|\mathbf{w}, \lambda) = p(\mathbf{x}|\mathbf{c}) p(\mathbf{c}|\mathbf{l}, \lambda_A) P(\mathbf{l}|\mathbf{w}, \lambda_L)$$

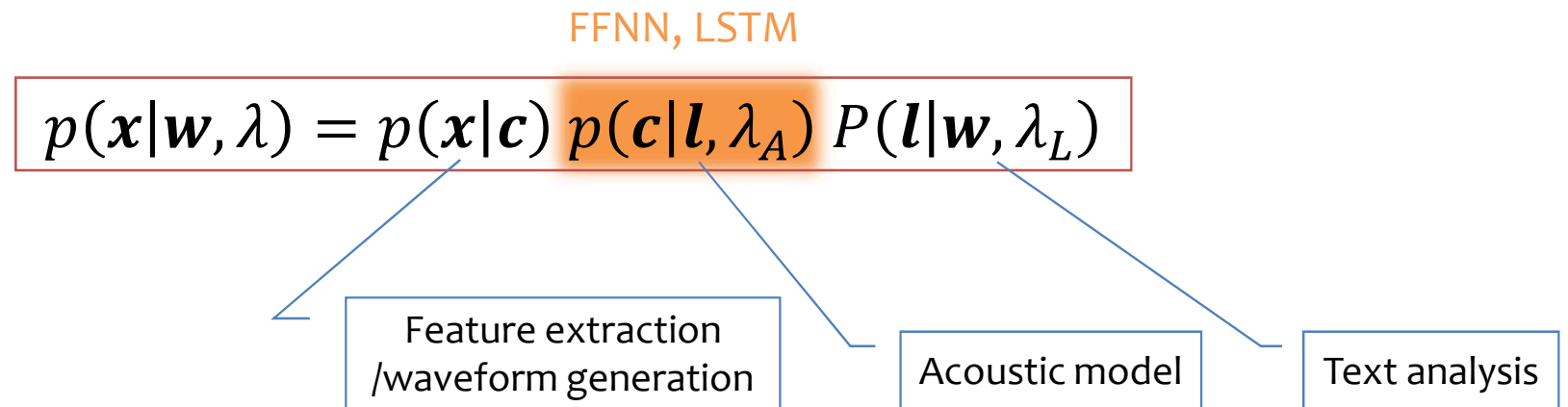

- Joint estimation of acoustic and excitation models [Maia '10]
- Log spectral distortion-version of MGE training [Wu '09]
- Factor analyzed trajectory HMM (STAVOCO) [Toda '08]
- Mel-cepstral analysis-integrated HMM [Nakamura '14]

Text analysis

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

Only from publications by the HTS working group

Deep neural network approaches (1/6)



- DNN-based speech synthesis [Zen '13]
- LSTM-based speech synthesis [Fan '14], etc.

DNN vs HMM

DNN

- Work for larger database?
- **Flat structure**
 - Easy to implement
 - Difficult to shouting troubles
- Often prior knowledge / model complexity is embedded in initialization and/or training process
- Suitable for parallel/distributed computation
- Optimization in continuous space

HMM (\cong regression tree)

- Can work for small database?
- **Plausible structure**
 - Difficult to implement
 - Easy to shouting troubles
- Prior knowledge / model complexity can be given in an explicit form (e.g., model structure)
- Unsuitable for parallel/distributed computation
- Optimization in discrete space

Deep neural network approaches (2/6)

Source
filter model FFNN, LSTM

$$p(\mathbf{x}|\mathbf{w}, \lambda) = p(\mathbf{x}|\mathbf{c}) p(\mathbf{c}|\mathbf{l}, \lambda_A) P(\mathbf{l}|\mathbf{w}, \lambda_L)$$

→

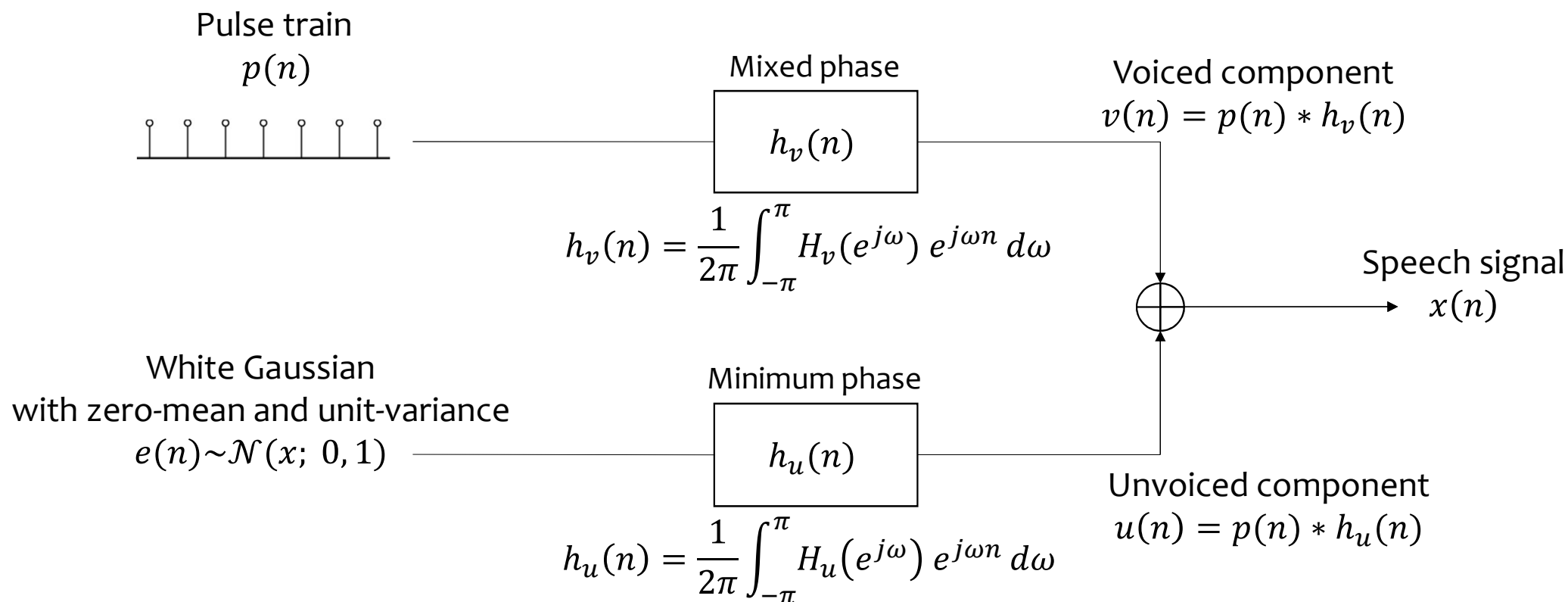
>>

1. Measuring likelihoods of speech waveform directly,
2. train a neural network
3. which models both voiced and unvoiced components.

t analysis

- Directly modeling speech waveforms by neural networks [Tokuda '15],
- Directly modeling voiced and unvoiced components by neural networks [Tokuda '16]

Speech signal model



Signal model for unvoiced+voiced sounds

Deep neural network approaches (3/6)

WaveNet, SampleRNN, WaveRNN, ...
(autoregressive structure)

$$p(\mathbf{x}|\mathbf{w}, \lambda) = p(\mathbf{x}|\mathbf{c}) p(\mathbf{c}|\mathbf{l}, \lambda_A) P(\mathbf{l}|\mathbf{w}, \lambda_L)$$

Feature extraction
/waveform generation

Acoustic model

Text analysis

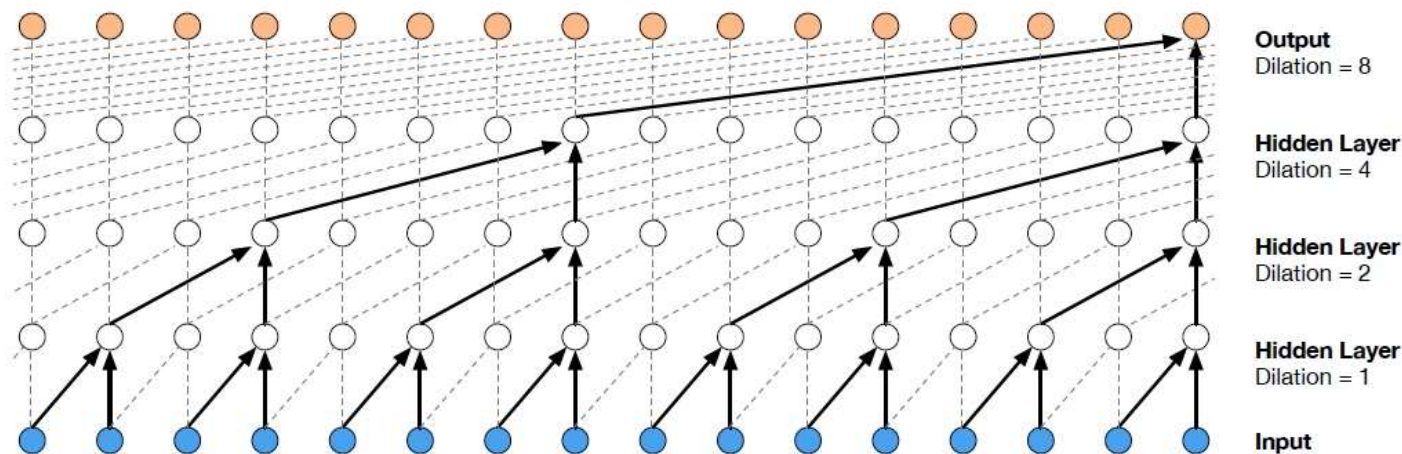
WaveNet

- Autoregressive generative model using convolutional NN
 - Directly modeling speech waveform

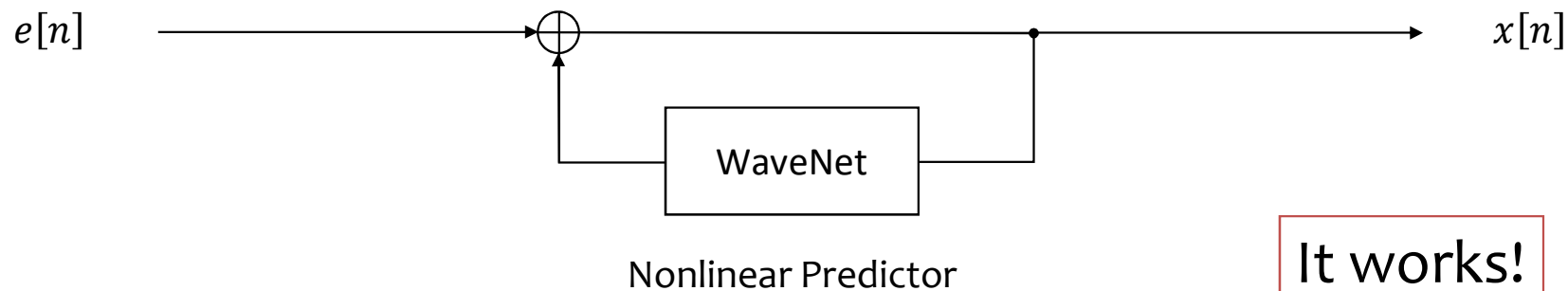
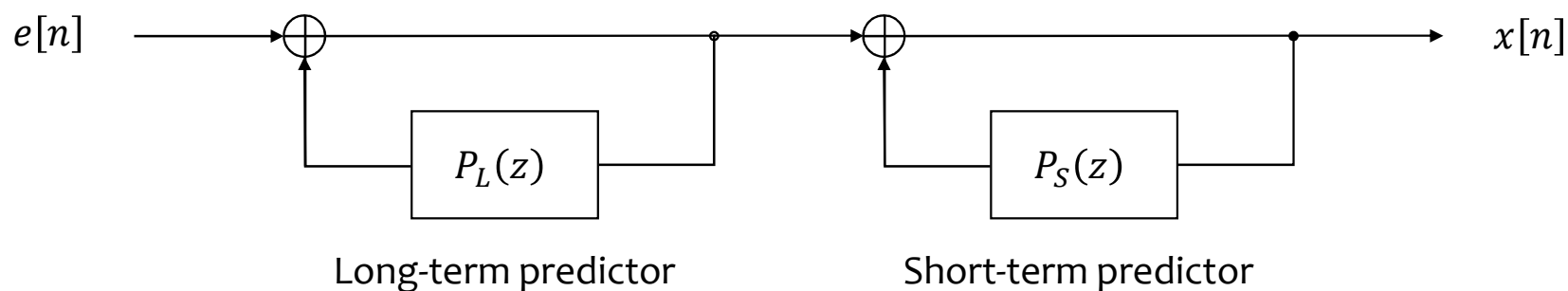
x : waveform h : acoustic and linguistic feature

$$p(x | h) = \prod_{n=0}^{N-1} \underbrace{p(x[n] | x[0], \dots, x[n-1], h)}_{\text{modeled by using CNN}}$$

- Dilated causal convolution



Speech signal generation model



It works!

Famous words in speech technology (1980s)

“Every time I fire a **linguist**,
the performance of the **speech recognizer** goes up”
by Frederick Jelinek

“Every time I fire a **speech technology researcher**,
the performance of the **speech synthesizer** goes up”
by ????? ?????

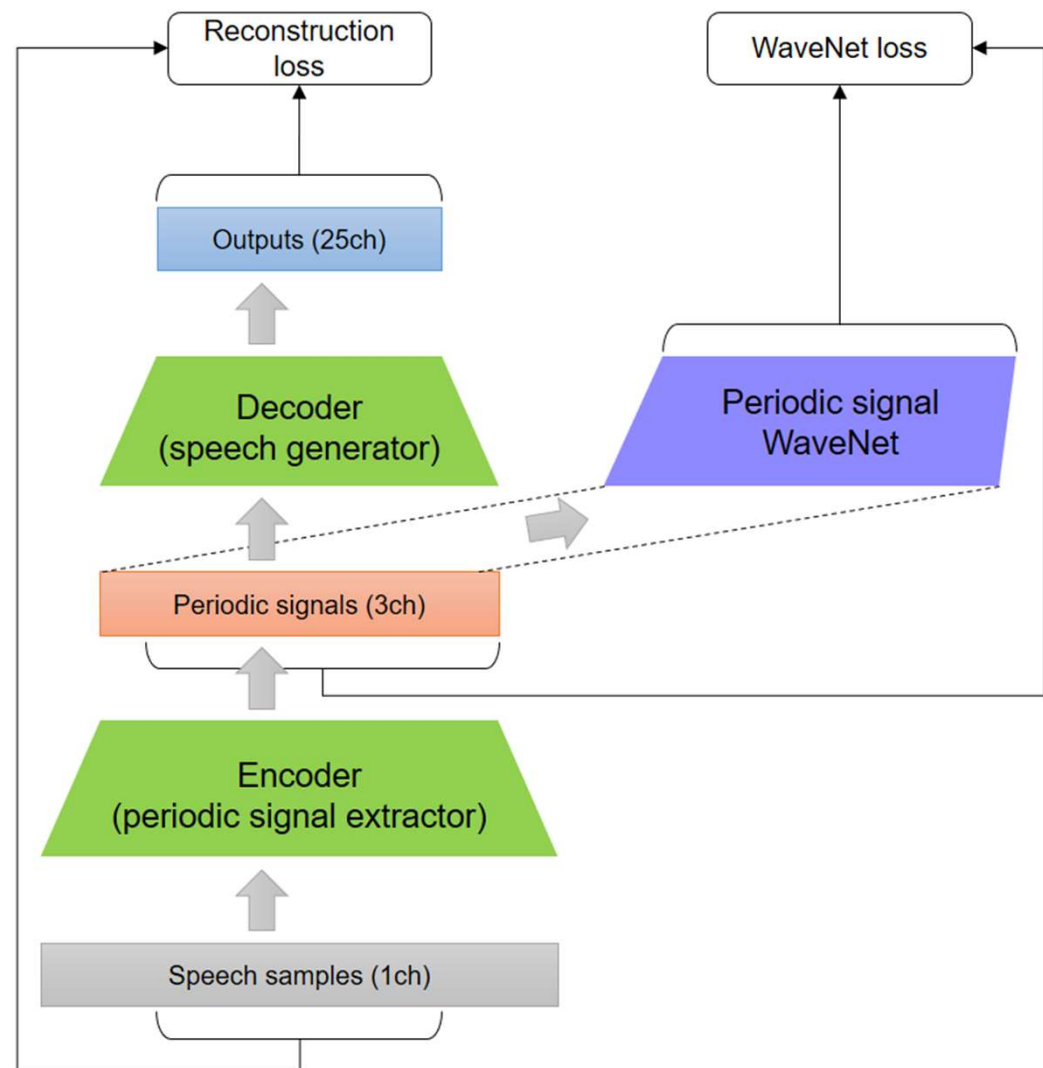
DNN variants for waveform modeling

- Autoregressive
 - WaveNet, SampleRNN, WaveRNN, ...
- Normalizing flow
 - WaveGlow, Parallel WaveNet, ClariNet, FloWaveNet, ...
- Combining with source filter model
 - LPCNet, ExcitNet, GlotNet, LP-WaveNet, ...
- Introducing signal processing technique
 - SubbandWaveNet, FFTNet, ...

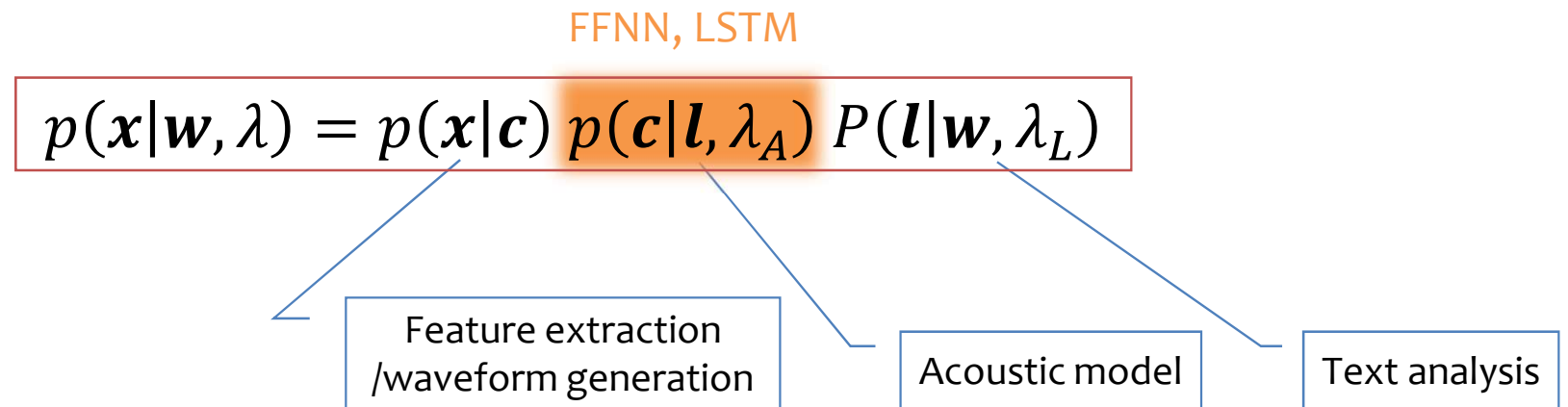


DNN vocoder with periodic excitation [Oura '19]

- Autoencoder-type structure extracts 3 dimensional periodic signal
- Decoder generates periodic components and stochastic components
- WaveNet gives a constraints on the intermediate variable



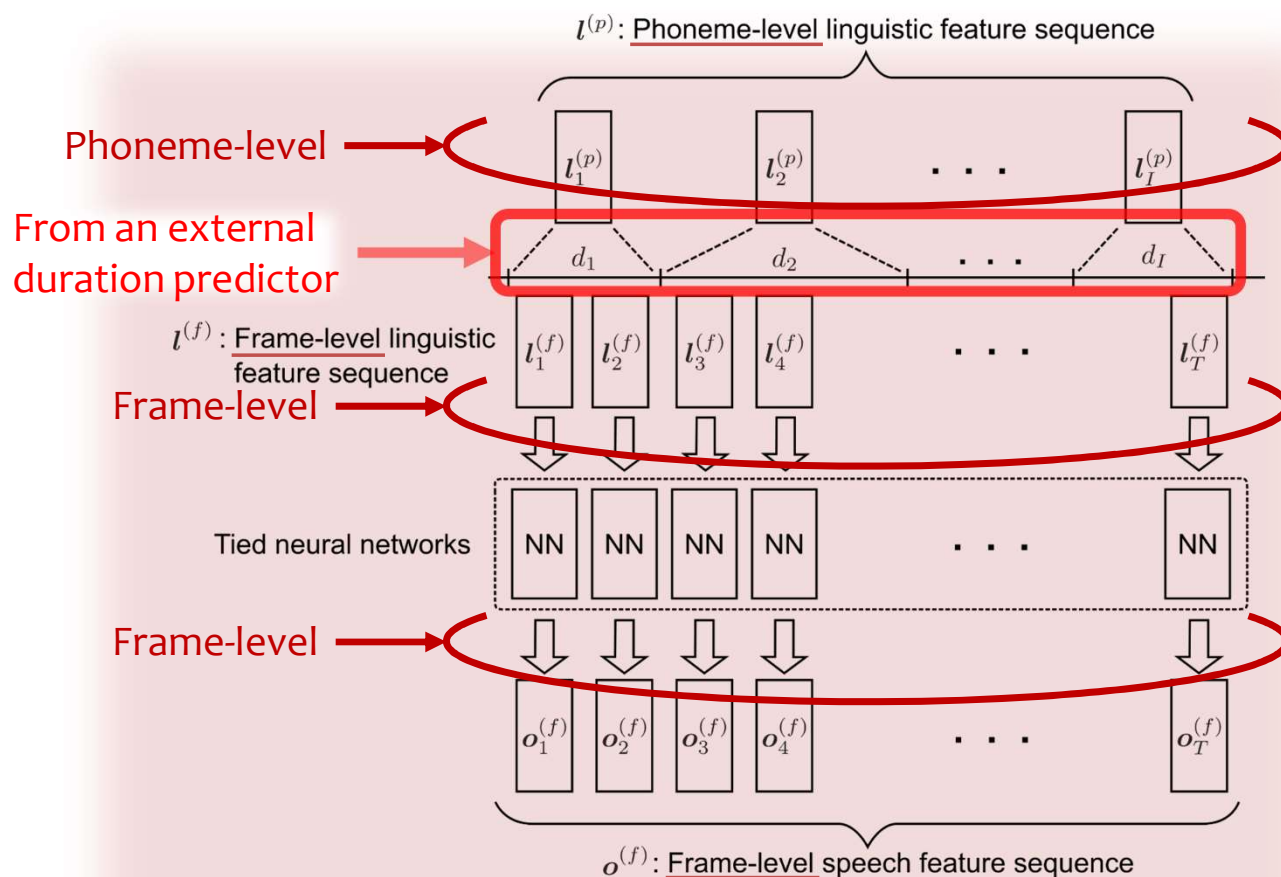
Deep neural network approaches (4/6)



- **HSMM**: duration model is included
- **FFNN, LSTM, WaveNet**: external duration predictor is required

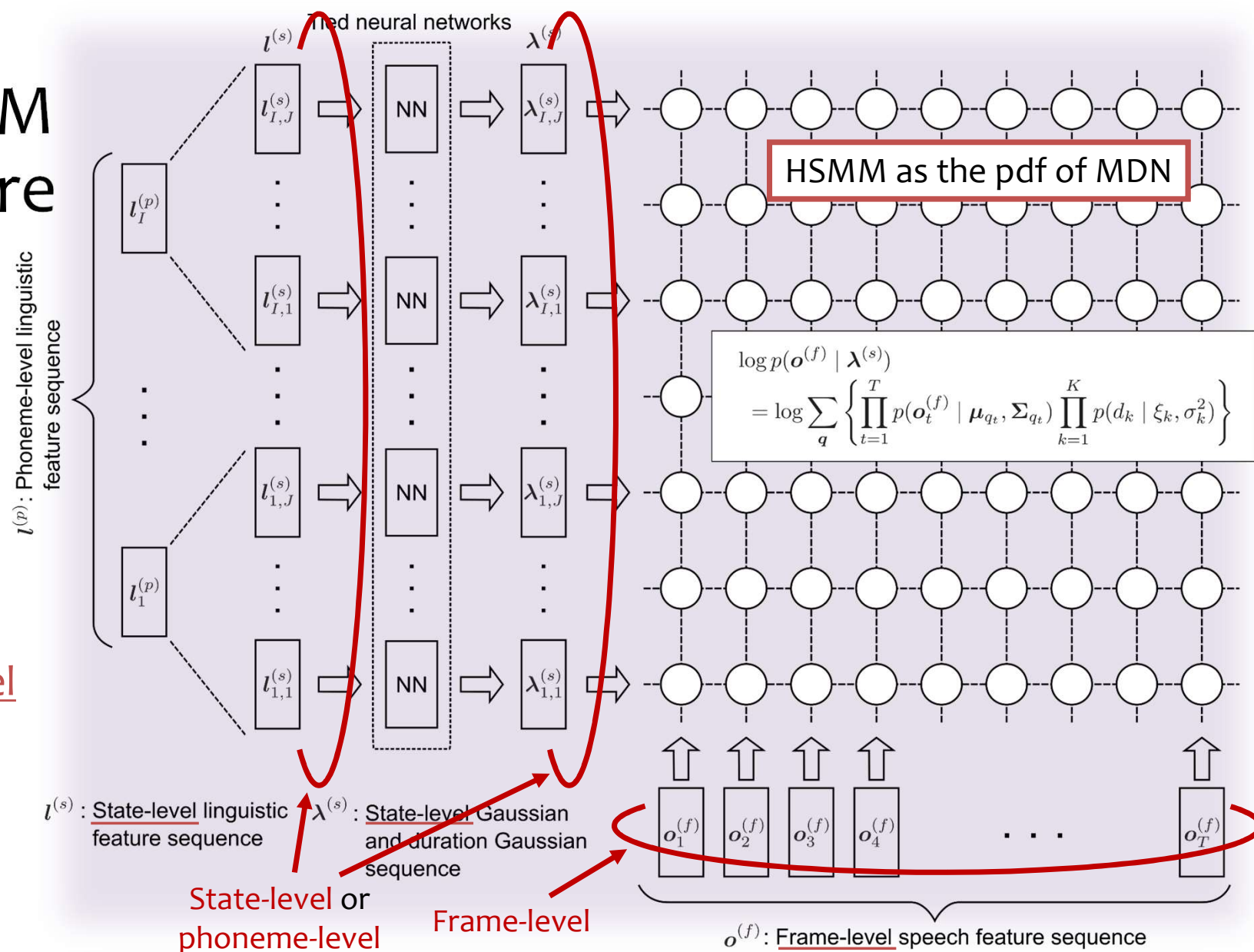
Frame-by-frame conversion

It needs an external duration predictor to determine phone durations



DNN-HSMM architecture

It runs at state-level
or phoneme-level
[Tokuda '16]



Deep neural network approaches (5/6)

WaveNet vocoder
(autoregressive structure)

Tacotron, Char2Wav, DeepVoice, ...
(attention mechanism)

$$p(\mathbf{x}|\mathbf{w}, \lambda) = p(\mathbf{x}|\mathbf{c}) p(\mathbf{c}|\mathbf{l}, \lambda_A) P(\mathbf{l}|\mathbf{w}, \lambda_L)$$

→

>>

Feature extraction
/ waveform generation

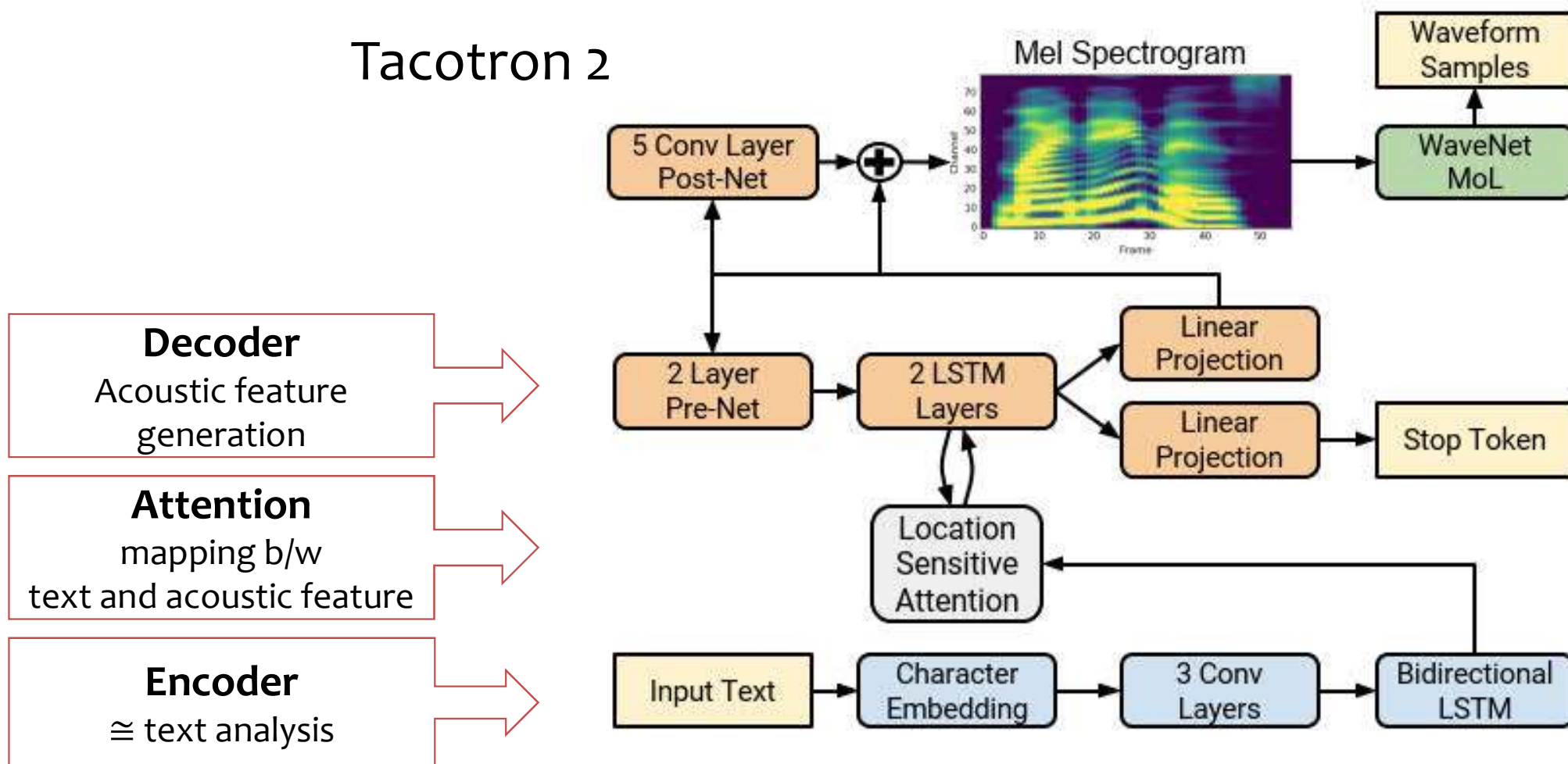
Acoustic model

Text analysis

Attention mechanism

[from [arXiv:1712.05884](https://arxiv.org/abs/1712.05884)]

Tacotron 2



Deep neural network approaches (6/6)

WaveNet vocoder
(autoregressive structure)

Attention

$$p(\mathbf{x}|\mathbf{w}, \lambda) = p(\mathbf{x}|\mathbf{c}) p(\mathbf{c}|\mathbf{l}, \lambda_A) P(\mathbf{l}|\mathbf{w}, \lambda_L)$$

Feature extraction
/waveform generation

Acoustic model

Text analysis



Controlling intermediate variables in the hierarchical structure

- Language
 - Japanese, English, Chinese, ...
- Dialect
- Pronunciation
- Pause
- Allophone
- Prosody
 - Accent, stress, tone, ...
- Speaking style, emotional expression
- Emphasis
- Nonverbal, paralinguistic information
- Voice characteristics
 - Male, female, child, adult, elderly
- Speech parameter
 - Fundamental frequency, volume, duration, aperiodic component, ...

High level

Text analysis

Acoustic model

Low level

Vocoding

Singing synthesis with CNN+WaveNet

WaveNet vocoder CNN

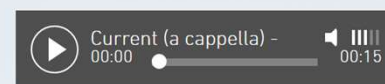
$$p(\mathbf{x}|\mathbf{w}, \lambda) = p(\mathbf{x}|\mathbf{c}) p(\mathbf{c}|\mathbf{l}, \lambda_A) P(\mathbf{l}|\mathbf{w}, \lambda_L)$$

Feature extraction
/ waveform generation

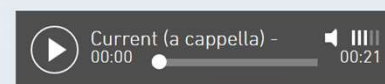
Acoustic model

Text analysis

[Japanese] Diamonds



[Japanese] 瞳 (Hitomi)



Other DNN techniques and architectures

- GAN
- VAE/VQ-VAE
- Transformer (self attention)

Outline

- Statistical formulation of speech synthesis
- HMM-based speech synthesis
- Deep neural networks
- **Evaluation / data & software tools**
- Other related topics

Blizzard Challenge

- Performance of TTS system depends on the database
- Difficult to compare techniques themselves



“Blizzard Challenge”

Evaluating corpus-based speech synthesis
on **common datasets** [Black '05]

Since 2005

Evaluation methodology

- Naturalness
 - Mean Opinion Score
- Speaker similarity
 - Degradation Mean Opinion Score
- Intelligibility (dictation of SUS, PCS, etc.)
 - Word accuracy

Not enough for spontaneous speech, audio book task, etc.

Section 1: Part 1 / 17

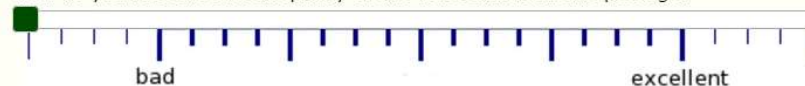
In this section, you will listen to a short passage from an children's audio book, and you will give your opinion about various aspects of the voice you just heard. You might like to imagine that you are choosing which of them to buy for a young child.



You will then choose a response for each question below. Your score will be represented by a slider. For example, the midpoint in the overall quality slider should be used to indicate that the quality is approximately half of the best possible quality.

Overall impression

How do you rate the overall quality of the voice that read this passage?



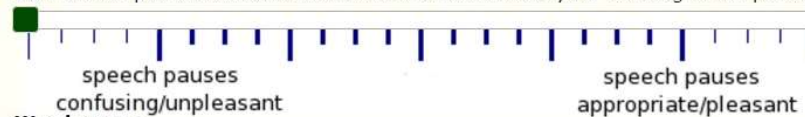
Pleasantness

How pleasant did you find the voice you just heard?



Speech pauses

How did the pauses between words and sentences affect your listening to the passage?



Word stress

What did you think of the way words in the passage were stressed?



Intonation

What did you think of the "melody" of the voice reading this passage?

Common datasets for speech synthesis

- ELRA <http://www.elra.info/>
- ELDA <http://www.elda.org/>
- LDC <https://www ldc.upenn.edu/>
- OpenSLR <http://www.openslr.org/>
- ARCTIC
- VCTK
- LibriTTS, ...

Not so many because it needs studio-quality recordings

Software tools

- ISCA SynSig <https://www.synsig.org/index.php/Software>
- ISCA SCOOT <https://www.isca-speech.org/iscaweb/index.php/scoot>

Software tools



Toolkit for building voice interaction systems



hts_engine

Speech synthesis engine



HTS

Training toolkit



Open JTalk

Japanese TTS system



SPTK

Speech signal processing toolkit



Sinsy

Singing synthesis system

Takashi Masuko, Noboru Miyazaki, Kazuhito Koishida, Takayoshi Yoshimura, Heiga Zen, Junichi Yamagishi, Keiichiro Oura, Akinobu Lee and others contributed

Outline

- Statistical formulation of speech synthesis
- HMM-based speech synthesis
- Deep neural networks
- Evaluation / data & software tools
- Other related topics
 - Text normalization
 - Voice conversion
 - Speech coding
 - Anti-spoofing
 - Physical simulation

Text normalization

- Text normalization is excluded from the end-to-end systems
- Still rule-based approach is the mainstream
- It would be included in the end-to-end process in the near future

Voice conversion

- Close relationship to speech synthesis
- DNN-approach has emerged also in voice conversion research
- Realtime application is essential
- Realtime (or low-latency) prosody conversion is a challenging problem

Speech coding

- WaveNet and other waveform modeling approaches seems to bring a revolution to speech coding.

- WaveNet based low rate speech coding [Kleijn '18]
- A Real-Time Wideband Neural Vocoder at 1.6 kb/s Using LPCNet [Valin '19]
- Low Bit-rate Speech Coding with VQ-VAE and WaveNet [Garbacea '19]
- High-quality speech coding with sample RNN [Klejsa '19]
- WaveNet-based zero-delay lossless speech coding [Yoshimura '18]
- Wavenet-based delay-free ADPCM Speech Coding [Yoshimura '19]

Imposture using speech synthesis

- Fear for spoofing with speech synthesis
 - On the security of HMM-based speaker verification systems against imposture using synthetic speech [Masuko '99]
- Detecting synthesized speech
 - A robust speaker verification system against imposture using an HMM-based speech synthesis system [Sato '01]
- ASVspoof 2015
 - [The First Automatic Speaker Verification Spoofing and Countermeasures Challenge](#)

Physical simulation vs Deep neural network

- In the future, techniques for measuring dynamics of vocal tract will be significantly progressed.
- Also, techniques for simulating speech production system will be progressed.



will it be possible to generate natural-sounding speech based on the physical simulation approach?

- Advantage: realistic constraints, lower dimensional representation
→ latent representation in DNN-based system?

Summary

Statistical approach to speech synthesis

- Now it has reached at the level that we cannot tell the difference between human and machine
- Still we have a lot of problems to be solved →
- More flexibility and controllability for realizing diversity of speech

Let us enjoy speech synthesis research!

Thank you!

>>

Speech synthesis in the future

- Spoken dialog system

Cross-lingual/multilingual

- **Let us enjoy speech synthesis research
for realizing diversity of speech!**

- Support for people with disabilities

Flexibility/diversity

- CALL

Thank you! voice language

- Content creation

Editor design

Common data

Special thanks

- **Supervisors:** Satoshi Imai, Tadashi Kitamura, Takao Kobayashi
- **Colleagues & students:** Takashi Masuko, Noboru Miyazaki, Takayoshi Yoshimura, Shinji Sako, Masatsune Tamura, Junichi Yamagishi, Tomoki Toda, Heiga Zen, Kazuhito Koishida, Tetsuya Yamada, Nobuaki Mizutani, Ryuta Terashima, Akinobu Lee, Keiichiro Oura, Keijiro Saino, Kenichi Nakamura, Yi-Jian Wu, Ling-Hui Chen, Shifeng Pan, Yoshihiko Nankaku, Ranniery Maia, Sayaka Shiota, Chiyomi Miyajima, Kei Hashimoto, Shinji Takaki, Kazuhiro Nakamura, Kei Sawada, Takenori Yoshimura, Daisuke Yamamoto, ...
- **Collaborators and advisors:** Junichi Takami, Naoto Iwahashi, Mike Schuster, Satoshi Nakamura, Frank Soong, Michael Picheny, Simon King, Steve Young, Mari Ostendorf, Alan Black, Alex Acero, Bill Byrne, Phil Woodland, Thomas Hain, Phil Garner, Masataka Goto, Shigeru Katagiri, Hideki Kenmochi, Kazuya Takeda, Tatsuya Kawahara, Sadaoki Furui, Seiichi Nakagawa, Keikichi Hirose, Tetsunori Kobayashi, Miko Kurimo, Shigeki Sagayama, Kiyohiro Shikano, Hisashi Kawai, Nobuyuki Nishizawa, Minoru Tsuzaki, Yoichi Yamashita, Nobuaki Minematsu, Mat Shannon, Mark Gales, Kai Yu, John Dines, ...

in random order. I am sorry but I may have missed many...

References

- K. Tokuda, T. Kobayashi, S. Shiimoto, S. Imai, “Adaptive filtering based on cepstral representation —adaptive cepstral analysis of speech,” ICASSP 1990a.
- K. Tokuda, T. Kobayashi, S. Imai, “Generalized cepstral analysis of speech — unified approach to LPC and cepstral method,” ISCSLP 1990b.
- T. Fukada, K. Tokuda, T. Kobayashi, S. Imai, “An adaptive algorithm for mel-cepstral analysis of speech,” ICSSP 1992.
- K. Tokuda, T. Kobayashi, T. Masuko, S. Imai, “Mel-generalized cepstral analysis —a unified approach to speech spectral estimation,” ICSLP 1994.
- K. Tokuda, T. Kobayashi, S. Imai, “Speech parameter generation from HMM using dynamic features,” ICASSP 1995a.
- K. Tokuda, T. Masuko, T. Yamada, T. Kobayashi, S. Imai, “An algorithm for speech parameter generation from continuous mixture HMMs with dynamic features,” EUROSPEECH 1995b.

Only from HTS group

- T. Masuko, K. Tokuda, T. Kobayashi, S. Imai, “HMM-based speech synthesis with various voice characteristics,” ASA/ASJ Joint Meeting 1996.
- T. Masuko, K. Tokuda, T. Kobayashi, S. Imai, “Voice characteristics conversion for HMM-based speech synthesis system,” ICASSP 1997.
- T. Yoshimura, T. Masuko, K. Tokuda, T. Kobayashi, T. Kitamura, “Speaker interpolation in HMM-based speech synthesis system,” EUROSPEECH 1997.
- T. Masuko, T. Kobayashi, M. Tamura, J. Masubuchi, K. Tokuda, “Text-to-visual speech synthesis based on parameter generation from HMM,” ICASSP 1998.
- M. Tamura, T. Masuko, K. Tokuda, T. Kobayashi, “Speaker adaptation for HMM-based speech synthesis system using MLLR,” SSW 1998.
- T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, T. Kitamura, “Simultaneous modeling of spectrum, pitch and duration in HMM-based speech synthesis,” EUROSPEECH 1999.
- K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, T. Kitamura, “Speech parameter generation algorithms for HMM-based speech synthesis,” ICASSP 2000.

- T. Masuko, T. Hitotsumatsu, K. Tokuda, T. Kobayashi, “On the security of HMM-based speaker verification systems against imposture using synthetic speech,” EUROSPEECH 1999.
- T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, T. Kitamura, “Simultaneous modeling of spectrum, pitch and duration in HMM-based speech synthesis,” EUROSPEECH 1999.
- T. Masuko, K. Tokuda, T. Kobayashi, “Imposture using synthetic speech against speaker verification based on spectrum and pitch,” ICSLP/INTERSPEECH 2000.
- S. Sako, K. Tokuda, T. Masuko, T. Kobayashi, T. Kitamura, “HMM-based text-to-audio-visual speech synthesis,” ICSLP/INTERSPEECH 2000.
- M. Tamura, T. Masuko, K. Tokuda, T. Kobayashi, “Text-to-speech synthesis with arbitrary speaker’s voice from average voice,” EUROSPEECH 2001.
- T. Satoh, T. Masuko, T. Kobayashi, K. Tokuda, “A robust speaker verification system against imposture using an HMM-based speech synthesis system,” EUROSPEECH 2001.
- K. Shichiri, A. Sawabe, K. Tokuda, T. Masuko, T. Kobayashi, T. Kitamura, “Eigenvoices for HMM-based speech synthesis,” ICSLP 2002.

- J. Yamagishi, T. Masuko, K. Tokuda, T. Kobayashi, “A training method for average voice model based on shared decision tree context clustering and speaker adaptive training,” ICASSP 2003.
- K. Tokuda, H. Zen, T. Kitamura, “Trajectory modeling based on HMMs with the explicit relationship between static and dynamic features,” EUROSPEECH 2003.
- H. Zen, K. Tokuda, T. Masuko, T. Kobayashi, T. Kitamura, “Hidden semi-Markov model based speech synthesis,” ICSLP 2004.
- T. Toda, A. Black, K. Tokuda, “Spectral conversion based on maximum likelihood estimation considering global variance of converted parameter,” ICASSP 2005.
- T. Toda, K. Tokuda, “Speech parameter generation algorithm considering global variance for HMM-based speech synthesis,” INTERSPEECH 2005.
- A. Black, K. Tokuda, “The Blizzard Challenge – 2005: Evaluating corpus-based speech synthesis on common datasets,” INTERSPEECH 2005.
- K. Saino, H. Zen, Y. Nankaku, A. Lee, K. Tokuda, “HMM-based singing voice synthesis system,” Interspeech 2006.

- R. Maia, T. Toda, H. Zen, Y. Nankaku, K. Tokuda, “An excitation model for HMM-based speech synthesis based on residual modeling,” SSW 2007.
- K. Oura, Y. Nankaku, T. Toda, K. Tokuda, R. Maia, S. Sakai, S. Nakamura, “Simultaneous Acoustic, Prosodic, and Phrasing Model Training for TTS Conversion Systems,” ISCSLP 2008.
- K. Hashimoto, H. Zen, Y. Nankaku, K. Tokuda, “A Bayesian approach to HMM-based speech synthesis,” ICASSP 2009.
- K. Kazumi, Y. Nankaku, K. Tokuda, “Factor analyzed voice models for HMM-based speech synthesis,” ICASSP 2010.
- K. Oura, A. Mase, T. Yamada, S. Muto, Y. Nankaku, K. Tokuda, “Recent Development of the HMM-based Singing Voice Synthesis System – Sinsy,” SSW 2010.
- K. Nakamura, K. Hashimoto, Y. Nankaku, K. Tokuda, “Integration of Acoustic Modeling and Mel-cepstral analysis for HMM-based Speech Synthesis,” ICASSP 2013.
- K. Tokuda, H. Zen, “Directly modeling voiced and unvoiced components in speech waveforms by neural networks,” ICASSP 2016.

- K. Tokuda, K. Hashimoto, K. Oura, Y. Nankaku, “Temporal modeling in neural network based statistical parametric speech synthesis,” SSW 2016.
- T. Yoshimura, K. Hashimoto, K. Oura, Y. Nankaku, K. Tokuda, “Mel-cepstrum-based quantization noise shaping applied to neural-network-based speech waveform synthesis,” IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 26, no. 7, pp. 1173-1180, July, 2018.
- J. Niwa, T. Yoshimura, K. Hashimoto, K. Oura, Y. Nankaku, K. Tokuda, “Statistical voice conversion based on WaveNet,” ICASSP 2018.
- T. Yoshimura, K. Hashimoto, K. Oura, Y. Nankaku, K. Tokuda, “WaveNet-based zero-delay lossless speech coding,” SLT 2018.
- T Yoshimura, K. Hashimoto, K. Oura, Y. Nankaku, K. Tokuda, “Speaker-dependent WaveNet-based delay-free ADPCM speech coding,” ICASSP 2019.
- S. Takaki, Y. Nankaku, K. Tokuda, “Spectral modeling with contextual additive structure for HMM-based speech synthesis,” SSW 2010.