Flexible speech synthesis based on hidden Markov models

Keiichi Tokuda Nagoya Institute of Technology

APSIPA ASC 2013, Kaohsiung November 1, 2013

Towards human-like talking machines

For realizing natural human-computer interaction, speech synthesis systems are required to have an ability to generate speech with:

- arbitrary speaker's voice
- speaking styles (e.g., reading, conversational)
- emotional expressions
 (e.g., happy, angry, sad)
- emphasis
- and so on



History

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

General unit-selection synthesis scheme



Overview of this talk

- 1. Introduction and background
- 2. Basic techniques in the system
- 3. Examples demonstrating its flexibility
- 4. Discussion and conclusion



HMM-based speech synthesis system



HMM-based speech synthesis system



Human speech production



air flow

Source-filter model



ML estimation of spectral parameter

Mel-cepstral representation of speech spectra



ML-estimation of mel-cepstrum

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

x : short segment of speech waveform (Gaussian process)

Synthesis Filter



Structure of MLSA filter

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



Waveform reconstruction



from the text to be synthesized \rightarrow <u>use of HMM</u>

HMM-based speech synthesis system



Hidden Markov model (HMM)



Structure of state output (observation) vector



Observation of F0



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



MSD-HMM for F0 modeling



Voiced / Unvoiced weights

Contextual factors

Phoneme

- {preceding, succeeding} two phonemes
- current phoneme

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

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]



Stream-dependent tree-based clustering (2)



HMM-based speech synthesis system



Speech parameter generation algorithm [Tokuda; '00]

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

$$P(o \mid \hat{l}, \hat{\lambda}) = \sum_{q} P(o \mid q, \hat{\lambda}) P(q \mid \hat{l}, \hat{\lambda})$$

$$\approx \max_{q} P(o \mid q, \hat{\lambda}) P(q \mid \hat{l}, \hat{\lambda})$$

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

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

Determination of state sequence (1/3)



Determine state sequence via determining state durations

Determination of state sequence

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

 $p_i(\cdot)$: state-duration distribution of i-th state d_i : state duration of i-th state K : # of states in a sentence HMM for \hat{l}

Gaussian

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

Speech parameter generation algorithm

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

$$P(o \mid \hat{l}, \hat{\lambda}) = \sum_{q} P(o \mid q, \hat{\lambda}) P(q \mid \hat{l}, \hat{\lambda})$$

$$\approx \max_{q} P(o \mid q, \hat{\lambda}) P(q \mid \hat{l}, \hat{\lambda})$$

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

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

Without dynamic feature



becomes a sequence of mean vectors
⇒ discontinuous outputs between states

Dynamic features

$$\Delta \boldsymbol{c}_{t} = \frac{\partial \boldsymbol{c}_{t}}{\partial t} \approx 0.5(\boldsymbol{c}_{t+1} - \boldsymbol{c}_{t-1})$$
$$\Delta^{2} \boldsymbol{c}_{t} = \frac{\partial^{2} \boldsymbol{c}_{t}}{\partial t^{2}} \approx \boldsymbol{c}_{t+1} - 2\boldsymbol{c}_{t} + \boldsymbol{c}_{t-1}$$



Integration of dynamic features

Relationship between speech parameter vectors & static feature vectors

$$\boldsymbol{o}_{t} = \begin{bmatrix} \boldsymbol{c}_{t}^{\mathsf{T}}, \Delta \boldsymbol{c}_{t}^{\mathsf{T}}, \Delta^{2} \boldsymbol{c}_{t}^{\mathsf{T}} \end{bmatrix}^{\mathsf{T}}$$



 \rightarrow_{29}

Solution for the problem (1/2)



we obtain

$$\boldsymbol{W}^{\mathsf{T}}\boldsymbol{\Sigma}_{\hat{q}}^{-1}\boldsymbol{W}\boldsymbol{c} = \boldsymbol{W}^{\mathsf{T}}\boldsymbol{\Sigma}_{\hat{q}}^{-1}\boldsymbol{\mu}_{\hat{q}},$$

where

$$\boldsymbol{c} = [\boldsymbol{c}_{1}^{\mathsf{T}}, \boldsymbol{c}_{2}^{\mathsf{T}}, \dots, \boldsymbol{c}_{T}^{\mathsf{T}}]^{\mathsf{T}}$$
$$\boldsymbol{\mu}_{\hat{q}} = [\boldsymbol{\mu}_{\hat{q}_{1}}^{\mathsf{T}}, \boldsymbol{\mu}_{\hat{q}_{2}}^{\mathsf{T}}, \dots, \boldsymbol{\mu}_{\hat{q}_{T}}^{\mathsf{T}}]^{\mathsf{T}}$$
$$\boldsymbol{\Sigma}_{\hat{q}} = [\boldsymbol{\Sigma}_{\hat{q}_{1}}^{\mathsf{T}}, \boldsymbol{\Sigma}_{\hat{q}_{2}}^{\mathsf{T}}, \dots, \boldsymbol{\Sigma}_{\hat{q}_{T}}^{\mathsf{T}}]^{\mathsf{T}}$$

Solution for the problem (2/2)



Generated speech parameter trajectory



Generated spectra



Spectra changing smoothly at state boundaries

Generated F0







Effect of dynamic features



HMM-based speech synthesis system



Overview of this talk

- 1. Introduction and background
- 2. Basic techniques in the system
- 3. Examples demonstrating its flexibility
- 4. Discussion and conclusion

Emotional speech synthesis

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

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 A's speaker-dependent system

Speaker interpolation (mixing voices)

Linear combination of two speaker-dependent models



Voice morphing

Two voices:



$A \diamond \diamond \diamond \diamond \diamond \diamond \diamond \diamond \diamond B last$

Four voices:



Interpolation of speaking styles



Eigenvoice (creating voices) [Shichiri; '02]



<u>Click here</u> for a demo

Multilingual speech synthesis

- Japanese 🛛 🐗 🐗
- American English 🛛 🐠 🐠 🐠 🐗
- Chinese (Mandarin) (by ATR) 4
- Brazilian Portuguese (by Nitech, and UFRJ) 4
- European Portuguese (by Nitech, Univ of Porto, and UFRJ) 4
- Slovenian (by Bostjan Vesnicer, University of Ljubljana, Slovenia)
- Swedish (by Anders Lundgren, KTH, Sweden) 4
- German (by University of Bonn, and Nitech)
- Korean (by Sang-Jin Kim, ETRI, Korea) \, 🐗
- Finish (by TKK, Finland) 🐗 🐗
- Baby English (by Univ of Edinburgh, UK)
- Polish, Slovak, Arabic, Farsi, Croatian, Polyglot, etc.

HMM-based singing synthesis



Flexibility in singing synthesis

• Sing a popular song



• Sing by a famous person's voice (from a TV program NHK "Science zero")



Rap singing



MMDAgent

A toolkit for building voice interaction systems

- Fully open-source toolkit with open interfaces
- Tightly integrated speech recognition/synthesis engines
- 3-D scene rendering fully compatible with CG tools
- HMM-based flexible and expressive speech synthesis (neutral, angry, bashful, happy, sad)







Summary

Statistical approach to speech synthesis

- Whole speech synthesis process is described in a unified statistical framework
- It can provide flexibility: various voices, emotional expressions, speaking styles, etc.
- Future work
- Still we have many problems should be solved:
 - Direct modeling of speech waveform
 - Importance of the database

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 vocaburary...," 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 guality," 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.

Keiichi Tokuda would like to thank HTS working group members, including Heiga Zen, Keiichiro Oura, Junichi Yamagichi, Tomoki Toda, Yoshihiko Nankaku, Kei Hahimoto, and Sayaka Shiota for their help. HTS Slides released by HTS Working Group <u>http://hts.sp.nitech.ac.jp/</u>

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 <u>http://creativecommons.org/</u> for details.



Copyright for some samples

Some samples (icon is <a>) are released under following license.

This voice is free for use for any purpose (commercial or otherwise) subject to the pretty light restrictions detailed below.

#######################################			
###		##	
###	Carnegie Mellon University	##	
###	Copyright (c) 2003	##	
###	All Rights Reserved.	##	
###		##	
###	Permission to use, copy, modify, and licence this software and its	##	
###	documentation for any purpose, is hereby granted without fee,	##	
###	subject to the following conditions:	##	
###	 The code must retain the above copyright notice, this list of 	##	
###	conditions and the following disclaimer.	##	
###	Any modifications must be clearly marked as such.	##	
###	Original authors' names are not deleted.	##	
###		##	
###	THE AUTHORS OF THIS WORK DISCLAIM ALL WARRANTIES WITH REGARD TO	##	
###	THIS SOFTWARE, INCLUDING ALL IMPLIED WARRANTIES OF MERCHANTABILITY	##	
###	AND FITNESS, IN NO EVENT SHALL THE AUTHORS BE LIABLE FOR ANY	##	
###	SPECIAL, INDIRECT OR CONSEQUENTIAL DAMAGES OR ANY DAMAGES	##	
###	WHATSOEVER RESULTING FROM LOSS OF USE, DATA OR PROFITS, WHETHER IN	##	
###	AN ACTION OF CONTRACT, NEGLIGENCE OR OTHER TORTIOUS ACTION,	##	
###	ARISING OUT OF OR IN CONNECTION WITH THE USE OR PERFORMANCE OF	##	
###	THIS SOFTWARE.	##	
###		##	

###		##	
###	See http://www.festvox.org/cmu_arctic/ for more details	##	
###		##	

55