

複合ウェーブレットモデルと F_0 パターン生成過程の確率モデルを用いたテキスト音声合成*

門脇健人¹, 北条伸克², 亀岡弘和^{1,2}
 (¹東大院・情報理工, ²NTT)

1 はじめに

本稿ではこれまで我々が提案してきたスペクトルパラメータ生成法と F_0 パターン生成法を併用したパラメトリックテキスト音声合成方式を提案する.

HMM (Hidden Markov Model) 音声合成方式 [1] では通常、メルケプストラム係数とそのデルタ量の特徴量とし、それらの系列を生成する HMM を仮定する。状態出力分布に正規分布または GMM (Gaussian Mixture Model) を仮定した場合、学習において状態出力分布の平均は同一状態に割り当てられたすべての時刻の観測特徴量の平均となる。メルケプストラムの平均をとる操作は対数スペクトルの平均をとる操作と等価であるため、各々が山と谷がはっきりしたスペクトルであってもそれらの平均をとると山と谷が平滑化されたスペクトルになってしまう。平滑化されたスペクトル包絡による合成音声は buzzy でこもった音になることが知られており、この問題の解決はパラメトリック音声合成における重要課題の一つである。[2] では、この問題に対し、従来のメルケプストラム系列を生成する GMM-HMM 系から脱却し、スペクトル系列を生成する複合ウェーブレットモデル (Composite Wavelet Model; CWM) [6, 7] と HMM の統合系を提案した。CWM はスペクトルのピークとパワーの両方に相当するパラメータを有していることが特色で、CWM-HMM 系の学習では各状態出力分布の平均はスペクトルのピークの周波数およびパワー方向に平均をとったものとなるため、従来のメルケプストラム GMM-HMM 系に比べ合成音声におけるスペクトル包絡の過剰平滑化が起きにくいのが特長である。

また、音声において音韻とともに韻律が果たす役割は重要で、HMM 音声合成方式では基本周波数 (F_0) パターンをいかに言語的に自然でかつ肉声らしいものにできるかも重要課題の一つである。従来の方式では、学習データの量によっては過学習を起こし不自然な F_0 パターンを生成することがあったが、言語的および物理的に不自然な F_0 パターンを生成することをできるだけ防ぎたいという動機から [3] では F_0 パターンの物理的生成過程を模した藤崎モデル [4] の確率モデル版 [8] を用いた韻律生成法を提案した。この方式では言語的かつ物理的に意味のあるパラメータが生成されるため、合成音声のイントネーションの話者性や発話スタイルを柔軟に加工できるという利点もある。

本稿では、当研究室で開発した上述の二つのモデルを併用したパラメトリック音声合成方式を提案し、それぞれの効果を検証する。

2 CWM-HMM 統合系

本章では、[2] で提案したスペクトル包絡系列の生成モデルについて概説する。

スペクトル包絡ピークの周波数とパワーをパラメータにもつ CWM により、その確率モデル化を行うことが可能である。CWM は、GMM によりスペクトル包絡を近似的に表現するモデルである。CWM によるス

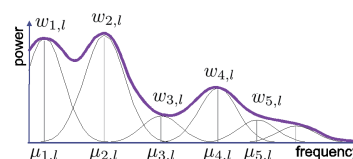


Fig. 1 CWM によるスペクトル表現

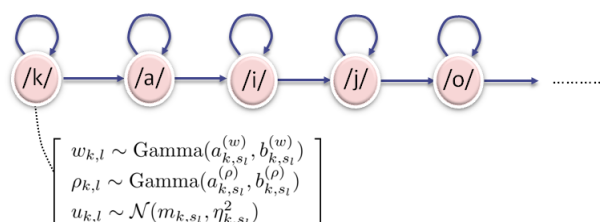


Fig. 2 CWM パラメータを出力する HMM の構成

ペクトル包絡 $f_{\omega,l}$ は、Fig. 1 のように、

$$f_{\omega,l} = \sum_{k=1}^K \frac{w_{k,l}}{\sqrt{2\pi}\sigma_{k,l}} \exp\left(-\frac{(\omega - \mu_{k,l})^2}{2\sigma_{k,l}^2}\right) \quad (1)$$

と表される。ただし、 ω, l は周波数と時刻のインデックス、 k は Gauss 関数のインデックスであり、 K は GMM の混合数である。また、 $\mu_{k,l}, \sigma_{k,l}^2, w_{k,l}$ は、それぞれ Gauss 分布関数を統計分布と見なした際の平均・分散・重みに対応し、各スペクトル包絡ピークの周波数・鋭さ・強度に対応するパラメータである。

以上の準備のもと、観測スペクトル系列が生成される過程を確率モデル化する。図 2 のような、各離散時刻 l ごとに平均 $\mu_{k,l}$ 、分散の逆数 $\rho_{k,l} := 1/\sigma_{k,l}^2$ 、重み $w_{k,l}$ の CWM パラメータを出力する HMM を考える。HMM の各状態は、言語ラベルの一状態を表しており、例えば Fig. 2 のようにそれぞれ一つの音素に対応させることもできるが、HTS [1] などの手法と同様に、音素状態に加え、前後の音素のやアクセント位置などの情報を用いたコンテクストラベルの一状態を対応させることも可能である。時刻 l における状態番号を s_l とし、本稿では状態出力分布 (各状態における CWM パラメータの生成確率) を、後述のパラメータ推定アルゴリズムの導出の便宜上の都合により、

$$P(\mu_{k,l}|s_l) = \mathcal{N}(\mu_{k,l}; m_{k,s_l}, \eta_{k,s_l}^2) \quad (2)$$

$$P(\rho_{k,l}|s_l) = \text{Gamma}(\rho_{k,l}; a_{k,s_l}^{(\rho)}, b_{k,s_l}^{(\rho)}) \quad (3)$$

$$P(w_{k,l}|s_l) = \text{Gamma}(w_{k,l}; a_{k,s_l}^{(w)}, b_{k,s_l}^{(w)}) \quad (4)$$

と仮定した。ここで、 $\mathcal{N}(x; m, \eta^2)$ は正規分布、 $\text{Gamma}(x; a, b)$ はガンマ分布

$$\text{Gamma}(x; a, b) = \frac{x^{a-1} \exp(-x/b)}{\Gamma(a) b^a} \quad (5)$$

*Text-to-speech synthesis combining CWM-HMM and probabilistic Fujisaki model. by KADOWAKI Kento (The Univ. of Tokyo), Hojo Nobukatsu (NTT), KAMEOKA Hirokazu (The Univ. of Tokyo/NTT)

である．また，上述のHMMにより生成されたCWMパラメータの系列 $\mu = \{\mu_k\}_{k,l}$, $\rho = \{\rho_k\}_{k,l}$, $w = \{w_k\}_{k,l}$ が与えられたとき，時刻 l において，観測スペクトル $y_{w,l}$ が生成される確率分布についても後述のパラメータ推定アルゴリズムの導出の便宜上の都合により，

$$P(y_{w,l}|\mu, \rho, w) = \text{Poisson}(y_{w,l}|f_{w,l}) \quad (6)$$

と定めた．ここで， $f_{w,l}$ は，cwmパラメータ系列 μ, ρ, w が与えられたとき，時刻 l のcwmパラメータを用いて(1)で表されるスペクトル包絡モデルであり， $\text{Poisson}(x; \lambda)$ はポアソン分布

$$\text{Poisson}(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad (7)$$

である．上記の生成モデルを定めることにより，次節のパラメータ推定アルゴリズムを導出する．

スペクトル包絡系列生成モデルパラメータの推定は，スペクトル包絡生成モデルのいわば逆問題である．これは，観測スペクトル系列 $Y = \{y_{w,l}\}_{w,l}$ が与えられたときに，スペクトル系列生成モデルパラメータ θ の事後確率 $P(\theta|Y) \propto P(Y|\theta)P(\theta)$ を最大化する問題として定式化される．推定すべきパラメータ θ は，HMMの出力状態列 ($s = \{s_l\}_l$)，HMMの各状態 i の出力分布パラメータ ($\theta = \{m_{k,i}, \eta_{k,i}, a_{k,i}^{(\sigma)}, b_{k,i}^{(\sigma)}, a_{k,i}^{(w)}, b_{k,i}^{(w)}\}_{k,i}$) と，CWMパラメータ系列 (μ, ρ, w) である．

θ の事後確率 $P(\theta|Y)$ を最大化する θ を求めることは難しいが，各変数について局所最適化を繰り返すことは可能である．このとき $P(\theta|Y)$ は，

$$\log P(\theta|Y) \stackrel{c}{=} \log P(Y|\theta) + \alpha \log P(\theta), \quad (8)$$

$$\log P(\theta) \stackrel{c}{=} \log P(s) + \log P(\rho|s, \theta) + \log P(w|s, \theta) + \log P(\mu|s, u, \theta) \quad (9)$$

と書ける．ここで，また， $\stackrel{c}{=}$ は定数部分を除いて一致することを意味する． $P(\theta|y)$ を最大化 (または $-P(\theta|y)$ を最小化) する θ を解析的に求めることは難しいが，補助関数法に基づいて局所最適化アルゴリズムを導くことができる．紙面の都合上詳細は省略するが，詳しくは [2] を参照されたい．

3 F_0 パターン生成過程の確率モデル

本章では， F_0 パターン生成過程の確率モデルに基づく韻律合成手法 [3] の概要を述べる．

藤崎モデル [4] は，甲状軟骨の二つの独立な運動 (平行移動運動と回転運動) に伴う声帯の伸びの長さの和が声帯の固有振動数の対数 ($\log F_0$) に比例する，という仮定をもとに，甲状軟骨の運動方程式を通して F_0 パターンの生成過程を表現したモデルである．藤崎モデルをベースにした F_0 パターンの生成過程の確率モデル [8, 9] は以下で定式化される． k を離散時刻のインデックスとし， $y_p[k]$, $u_p[k]$, $y_a[k]$, $u_a[k]$ をそれぞれフレーズ成分 $y_p(t)$ ，フレーズ指令 $u_p(t)$ ，アクセント成分 $y_a(t)$ ，アクセント指令 $u_a(t)$ の離散時間表現として，観測 F_0 パターンの対数値 $y[k]$ を次のように表現する．

$$y[k] | u_p[k], u_a[k] \sim \mathcal{N}(x[k], v_n^2[k]), \quad (10)$$

$$x[k] = G_p[k] * u_p[k] + G_a[k] * u_a[k] + u_b. \quad (11)$$

ここで $v_n^2[k]$ は時刻 k における観測 F_0 パターンの“不確かさ”を表すために導入した変数であり，これにより全時刻で正しい F_0 の値が観測できるとは限らないという問題をノイズとして統一的に扱うことを可能にした．次に， $u_p[k]$ と $u_a[k]$ は，それぞれイン

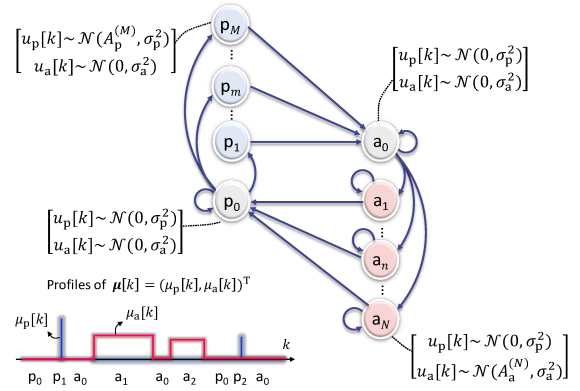


Fig. 3 指令列生成 HMM の状態遷移モデル [3] .

パルス列状および矩形パルス列状の指令列関数であり，各パルスが同時刻には生起しない，という制約を満たす必要がある．[8, 9] では，両指令列関数のペア $o[k] = (u_p[k], u_a[k])^T$ を以下に示す HMM (以後，指令列生成 HMM) の出力系列と見なそうというアイデアにより，上述の制約を満たした指令列関数の確率モデルが提案されている．

[3] でも述べたように，藤崎モデルの指令列は言語情報との関連が深いので，フレーズ指令と呼気段落，アクセント指令とアクセント句が 1 対 1 に対応した状態へと拡張することによって，各指令の強度をコンテキストに依存するパラメータとして扱うことで，コンテキストラベルを用いた統計学習により指令強度とコンテキスト情報を紐付けた．

上述の確率モデルを，コンテキスト依存型のモデルとするため Fig. 3 に示すような HMM の状態遷移を考える．このようなモデルは新たに以下の HMM で表現できる．

出力系列: $o[k] = (u_p[k], u_a[k])^T$ ($k = 1, \dots, K$)
状態集合: $S = \{p_0, \dots, p_M, a_0, \dots, a_N\}$
状態系列: $s = \{s_k \in S k = 1, \dots, K\}$
出力分布: $P(o[k] s_k = i) = \mathcal{N}(c_i[k], \Upsilon)$
$c_i[k] = \begin{cases} (0, 0)^T & (i \in p_0, a_0) \\ (A_p^{(m)}, 0)^T & (i \in p_m) \\ (0, A_a^{(n)})^T & (i \in a_n) \end{cases} \quad \Upsilon = \begin{bmatrix} v_{p,i}^2 & 0 \\ 0 & v_{a,i}^2 \end{bmatrix}$
遷移確率: $\phi_{i',i} = \log P(s_k = i' s_{k-1} = i)$

以上の指令列生成 HMM において，フレーズ指令，アクセント指令の生起時刻がそれぞれ呼気段落とアクセント句に対応する様な制約を加えるために，Left-to-Right 型の HMM をコンテキストラベルをもとに設計することによって，制約を考慮しつつ，パラメータ学習において指令列の強度と状態系列を未知パラメータとして観測 F_0 パターンから推定することを考える [15]．ここで，以下の文字をまとめて，

$$y = \{y[k]\}_{k=1}^K, s = \{s_k\}_{k=1}^K,$$

$$o = \{(u_p[k], u_a[k])^T\}_{k=1}^K,$$

$$\theta = \{\{A_p^{(m)}\}_{m=1}^M, \{A_a^{(n)}\}_{n=1}^N\},$$

$$\text{及び } \mu_p = (\mu_p[1], \dots, \mu_p[K])^T, \mu_a = (\mu_a[1], \dots, \mu_a[K])^T,$$

$$y = (y[1], \dots, y[K])^T,$$

と表記する．簡単のため $\phi_{i',i}$, μ_b , $v_{p,i}^2$, $v_{a,i}^2$, v_b^2 , $v_n^2[k]$, α , β は定数と仮定すると，指令列生成 HMM の状態系列 s と状態出力分布パラメータ θ が与えられた下

Table 1 従来法と提案法によって生成された F_0 パターンと観測 F_0 パターンとの RMSE と相関係数を比較した表。括弧内はその標準偏差を表している。RMSE の低い方、相関係数の高い方を太字で表示した。

手法	RMSE	Correlation
提案手法	0.100 (± 0.023)	0.720 (± 0.458) × 10 ⁴
従来手法 [1]	0.122(± 0.033)	0.288(± 0.188) × 10 ⁴

5.2 主観評価

本節では、提案した音声合成システムが高品質な合成音声を生成できるかどうかを確かめるために行った主観評価実験について述べる。より細かく性能を評価するため、以下の四手法を評価対象とした。

Table 2 主観評価の対象とする四手法。ただし、hts は従来手法 [1]、cwm は cwm-hmm 統合系モデルの略で 2 節で述べた手法、pfm は probabilistic fujisaki model の略で 3 節で述べた手法である。

手法名	スペクトル包絡系列	F_0 パターン
A	hts	hts
B	cwm	hts
C	hts	pfm
D	cwm	pfm

以上の四手法に対してプリファレンステスト (AB テスト) を実施して、合成音声の品質を評価した。ここで、品質を評価する項目に関しては、JEITA 規格 [14] における評価項目から選んだ四つの項目、具体的には (a) 音声の明瞭性：音質の良さ、(b) 音声の時間要因：日本語らしいリズム感、(c) 音声の抑揚：抑揚・アクセントの自然性、流暢さ、(d) 音声の総合的な良さ：全体的な良し悪し、を用いた。この四手法に関して、受聴者は 4 手法全ての組み合わせの音声を聴き、各項目について良いと思った方を選択させる。受聴者は男性 5 名とする。実験は音声は HTS2.1 のデモスクリプトの j セット中の 5 文に対して行った。各項目に関するプリファレンススコアを以下の表に示す。表の各行が他の各手法と比較した場合のプリファレンススコアを表している。また、二項検定により 5% 水準で統計的に有意にスコアが高いと評価されたスコアは太字で示してある。

Table.3 から分かるとおり、提案音声合成手法は項目 (a) ~ (d) に関して、他の全ての手法と比較して有意に高いスコアを示した。また、手法 B や C は手法 A と比較して統計的に有意に高いとは言えなかったが、手法 D では他の全ての手法と比べて有意に高いスコアを示していることから、音韻情報と韻律情報の相乗的な効果により音声の自然性が飛躍的に良くなるという可能性が考えられる。

6 まとめ

本稿では、より高品質なテキスト音声合成システムを実現するため、CWM-HMM 統合モデルと F_0 パターン生成過程の確率モデルによる音声合成システムを提案し、その概要について述べた。性能評価においては、 F_0 パターン生成手法について客観評価基準を用いて評価し、その妥当性を示した。さらに、主観評価実験によって、提案システムの効果を確認した。今後は、更に大規模な主観評価実験を行う予定である。

謝辞 本研究は JSPS 科研費 26280060, 26730100 の助成を受けたものです。

Table 3 項目 (a) ~ (d) に関する各手法におけるプリファレンススコア。太字で示されているスコアは二項検定により 5% 水準で統計的に有意に高いと評価されていることを表す。

(a) 音質の良さ				
	A	B	C	D
A	-	0.500	0.400	0.033
B	0.500	-	0.400	0.167
C	0.600	0.600	-	0.200
D	0.967	0.833	0.800	-

(b) 日本語らしいリズム感				
	A	B	C	D
A	-	0.467	0.433	0.300
B	0.533	-	0.400	0.233
C	0.567	0.600	-	0.300
D	0.700	0.767	0.700	-

(c) アクセントの自然性				
	A	B	C	D
A	-	0.600	0.367	0.400
B	0.400	-	0.333	0.200
C	0.633	0.667	-	0.433
D	0.600	0.800	0.567	-

(d) 全体的な良し悪し				
	A	B	C	D
A	-	0.500	0.400	0.033
B	0.500	-	0.400	0.167
C	0.600	0.600	-	0.200
D	0.967	0.833	0.800	-

参考文献

- [1] K. Tokuda *et al.*, in Proc. ICASSP, vol. 3, pp. 1315–1318, 2000.
- [2] N. Hojo *et al.*, In *8th ISCA Workshop on Speech Synthesis*, pp.129–34, 2013.
- [3] K. Kadowaki *et al.*, in Proc. *The 15th Annual Conference of the International Speech Communication Association (Interspeech 2014)*, pp. 2322–2326, Sep. 2014.
- [4] H. Fujisaki, Raven Press, 1988.
- [5] 門脇他, 日本音響学会春季研究発表会講演集, 3-6-17, pp. 361–364, Mar. 2014.
- [6] 槐他, 信学技報, Vol. 105, No. 372, pp. 1–6, 2005.
- [7] H. Kameoka *et al.*, IEEE Transactions on Audio, Speech and Language Processing, Vol. 18, No. 6, pp. 1507–1516, Aug. 2010.
- [8] H. Kameoka *et al.*, in Proc. SAPA, pp. 43–48, 2010.
- [9] K. Yoshizato *et al.*, in Proc. Interspeech, 2012.
- [10] T. Mausko *et al.*, IEIC Technical Report, vol. 101, no. 323, pp. 41–42, 2001.
- [11] <http://hts.sp.nitech.ac.jp/>
- [12] H. Kawahara *et al.*, Speech Communication, vol. 27, no. 3, pp. 187–207, 1999.
- [13] T. Yoshimura *et al.*, in Proc. of Eurospeech, pp. 2347–2350, 1999.
- [14] <http://www.jeita.or.jp/>
- [15] 門脇 健人, 亀岡 弘和, 日本音響学会 2014 年秋季研究発表会講演集, 3-7-3, pp. 261–264, Sep. 2014.