

# 隠れマルコフモデルに基づいた歌声合成システム

酒 向 慎 司<sup>†</sup> 宮 島 千 代 美<sup>†</sup>  
 徳 田 恵 一<sup>†</sup> 北 村 正<sup>†</sup>

隠れマルコフモデルに基づく音声合成方式を歌声合成に拡張することにより構築した歌声合成システムについて述べる．本システムでは，歌い手の声の質と基本周波数パターンに関する特徴をモデル化するため，スペクトルと基本周波数パターンを HMM により同時にモデル化している．特に，自然な歌声を合成するうえで重要な要素となる音符の音階や音長の基本周波数パターンへの影響を精度良くモデル化するため，楽譜から得られる音階と音長を考慮したコンテキスト依存モデルを構築している．これらのモデルに対して決定木によるコンテキストクラスタリングを行うことで，未知の楽曲からの歌声合成が可能となっている．実験から，歌い手の特徴を再現し歌声の合成が可能であることを示す．

## A Singing Voice Synthesis System Based on Hidden Markov Model

SHINJI SAKO,<sup>†</sup> CHIYOMI MIYAJIMA,<sup>†</sup> KEIICHI TOKUDA<sup>†</sup>  
 and TADASHI KITAMURA<sup>†</sup>

We describe a singing voice synthesis system by applying HMM-based speech synthesis technique. In this system, a sequence of spectrum and F0 are modeled simultaneously in a unified framework of HMM, and context dependent HMMs are constructed by taking account of contextual factors that affects singing voice. In addition, the distributions for spectral and F0 parameter are clustered independently by using a decision-tree based context clustering technique. Synthetic singing voice is generated from HMMs themselves by using parameter generation algorithm. In the experiments, we confirmed that smooth and natural-sounding singing voice is synthesised. It is also maintains the characteristics and personality of the donor of the singing voice data for HMM training.

### 1. はじめに

今日，様々なテキスト音声合成システムが開発され，人々の身近なところで利用されつつある．また，品質向上にとどまらず，個人性や感情といった豊かな表現を可能とするための研究が各所で進められており，コンピュータとの対話手段のほかにも様々な用途が期待されている．たとえば，コンピュータによる自然な歌声の合成は，エンタテインメントやアミューズメント分野への応用を考えることができ，歌声を合成する試みは，これまでもいくつか提案されている<sup>1),2)</sup>．これからの歌声合成システムには，合成された音声の品質だけでなく，音楽性の豊かな表現能力が求められているが，現在の音声合成技術のレベルから考えると，まだ解決しなくてはならない問題が多く残されている．本論文では，音楽表現に表れる個人性に着目し，個

人の音楽表現を再現できる歌声合成システムの構築について検討する．入力された楽譜に従って，特定の歌い手の声質，スタイルの歌声を合成することで，テレビゲーム，カラオケなどのエンタテインメント，アミューズメントシステム，または玩具などに個性のある新たな表現を加えることができるほか，楽曲作成において曲のデモンストレーションを計算機上で容易に試行できるなどの応用を考えることができる．

このような歌声合成システムを構築するために，テキスト音声合成を利用することが考えられる．これまでに提案されてきた音声合成システムの多くは単位選択という方式に分類される．これは音素や音韻といった音声単位ごとに分類した波形データを，合成したいテキストに従ってつなぎ合わせることで音声合成する手法である．発声された音声波形を利用できるため，クリアな合成音を得やすいというメリットがある一方，接続部分の歪みが生じやすい，多様な声質や発話スタイルなどを得ようとすると，膨大な波形データを必要とするなどの問題がある．それに対して，隠れマルコ

<sup>†</sup> 名古屋工業大学大学院工学研究科  
 Department of Computer Science and Engineering,  
 Nagoya Institute of Technology

フモデル (Hidden Markov model: HMM) に基づいた音声合成手法<sup>3)</sup>では、これらの問題を解決するため、音声認識の分野で広く利用されている HMM を利用し、HMM 自身から音声パラメータを生成する。この手法の特長として、動的特徴量を考慮したパラメータ生成アルゴリズム<sup>4)</sup>によって、滑らかに変化する音声パラメータが得られるほか、モデルパラメータの変換により、別の話者への適応や、多様な声質や感情を表現した音声の合成に柔軟に対応できるなどの点があげられる<sup>5)~7)</sup>。

本論文では、我々がこれまでに提案してきた HMM に基づいた音声合成手法を拡張した歌声合成手法について検討し、歌い手の話者性や歌い方の特徴を再現可能な歌声合成システムを構築することを目的とする。本システムの大きな特徴は、システムのすべてのモデルパラメータを学習データ提供者の歌声により自動学習する点にある。音の高さや長さは、合成時に楽譜から一意に定めることもできるが、そこから合成される歌声は単調で機械的なものになり、歌声としての個性や魅力に欠けるものである。実際の歌声は、歌唱法や歌い手独自の特徴によって様々な形で表現されるものであり、音楽としての重要な要素となっている。それらの多彩な表現の特徴は、主に声の高さや時間的な構造などに表れていると考えられる。そこで、自然の歌声にあるような音の高さの変化を再現するため、声質を表すスペクトル情報と高さを表す基本周波数を、可変次元に対応した多空間上の確率分布に基づく HMM (Multi-Space probability distribution HMM; MSD-HMM)<sup>8)</sup> を用いてモデル化している。

さらに、より精密なモデル化を行うために、コンテキスト依存モデルを学習する。歌声は、通常の会話やテキストの読み上げなどの場合と比較して、発声する音の高さや時間的な長さ、声の強弱などの変動の様子が大きく異なることから、歌声に特化したモデル化手法が必要となる。歌い手は、歌詞、音高、音長、さらには、発想記号や表情記号など、楽譜から得られる様々な情報を基に、多様な歌唱表現を行う。本論文では、歌声データベースのサイズの問題や、歌声モデル自動構築に関する初めての試みであることなどを考慮し、歌詞のほか、音高と音長をコンテキストとして考え、前後の環境を考慮したそれらの組合せについてモデルを分類したコンテキスト依存モデルを用いることとする。続いて、前後の環境を考慮したそれらの組合せについてモデルを分類したコンテキスト依存モデルとし、これらのモデルに決定木によるコンテキストクラスタリング<sup>9)</sup>を適用することにより、未知の楽譜に対して

も自然な歌声の合成を可能としている。

本システムはボコーダベースのシステムであることから、生成される音質には一定の限界があり、実際の歌声と同等の品質が得られるものではない。しかし、学習データに基づいて得られた歌声モデルからは、楽譜上では表現できない歌い手の持つ様々な特徴を備え、合成時にそれらを再現することが可能となる。これはロボットや玩具などへ新たな個性を付加することができるという点で様々な応用が期待できる。実験では、童謡など 60 曲を収録して構築した歌声データベースから歌声モデルを学習し、本論文では、特に未知の楽曲からでも、学習データの歌い手の特徴を再現した歌声の合成が可能であることを示す。

以下、本論文は次のように構成されている。2 章では HMM に基づいた歌声合成システム、3 章でデータベースの収録と整備、4 章で実験および合成された歌声の評価と考察を述べ、最後に 5 章でまとめる。

## 2. HMM 歌声合成システム

本研究で提案する歌声合成システムの概略を図 1 に示す。本システムは、大きく、学習部と、歌声合成部の 2 つに分けられる。学習部では、初期モデルを基に楽譜情報と歌の波形データからなる歌声データベースを用いて歌声モデルを学習し、合成部では、合成したい歌の楽譜情報を入力として、学習部で得られた歌声モデルから歌声を合成する。なお、各部で利用される楽譜データとして、MIDI<sup>10)</sup> を利用している。

### 2.1 学習部

学習部では、音声合成用のモデルとして、スペクトルパラメータ、基本周波数、および継続長を HMM によって音素単位でモデル化する。声質を表すパラメータには様々な分析手法が提案されており、これらのスペクトルパラメータは、連続 HMM によってモデル化することができる。一方、音高を表すパラメータで

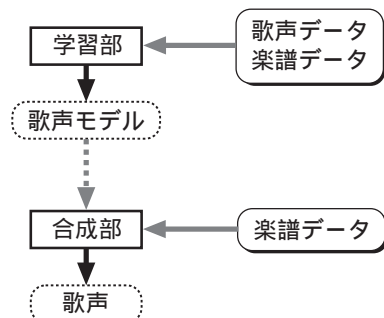


図 1 歌声合成システムの概略

Fig. 1 Block diagram of singing voice synthesis system.

ある基本周波数は有声区間では連続値をとり、無声区間では値を持たない可変次元の時間系列信号であるため、通常の連続 HMM や離散 HMM で直接モデル化することはできない。そこで、可変次元に対応した多空間上の確率分布に基づく MSD-HMM<sup>8)</sup> を用いて、スペクトルパラメータとしてメルケプストラム<sup>11)</sup> を多次元ガウス分布、基本周波数の有声音を 1 次元空間、無声音を 0 次元空間のガウス分布として単一の枠組みの中で同時にモデル化する。メルケプストラムは音声認識でよく用いられる特徴量であるが、歌声の表現に適したものであるかどうかは検討の余地がある。しかし、この問題に関して検討する点は多岐にわたることから、本論文では、これまでの HMM 音声合成で用いられてきたメルケプストラムを特徴パラメータとして使用することとする。

また、歌声に表れる声の特徴は、様々な要因によって影響を受け変動していると考えられる。たとえば、同じ音階の声であっても、広い範囲では楽曲のジャンルやテンポ、局所的には前後の歌詞や音階などによって、異なる特徴を持っていると考えられる。テキスト音声合成においても、テキストから得られる言語的な情報が音声パラメータに影響を与えていると考え、それらの要因をコンテキストと呼び、コンテキストを考慮したモデル化が行われている。

コンテキストに依存したモデル化を行うことで、精度の高い歌声モデルを得ることができるが、コンテキストの種類に応じてその組合せの数も莫大となってしまう。また、すべてのコンテキストの組合せに対応したモデルについて、十分な学習を行うためには、あらゆるパターンを網羅したデータベースが必要となってしまうため、現実的ではない。

この問題に対する優れた解決法として、コンテキストクラスタリングによってモデル間でパラメータを共有させる手法がある<sup>12)</sup>。これは、二分木を用いて、モデルの集合を木構造に分割することで、類似したコンテキストの組合せごとにモデルパラメータをクラスタリングする手法である。木の各ノードには、コンテキストを二分する質問があり、各リーフノードには、特定のモデルに相当するモデルパラメータがある。任意のコンテキストの組合せは、ノードにある質問に沿って木をたどることで、何らかのリーフノードに到達でき、該当するモデルを選択することができる。

たとえば、図 2 のように、隣り合う前後の音素をコンテキストとして、音素の三つ組みで表されたモデル集合をクラスタリングした場合、ある音素の組合せは、各ノードにある音素の三つ組みを分類する質問に

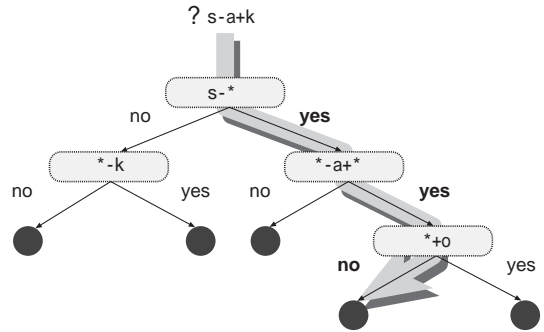


図 2 決定木によるクラスタリング  
Fig. 2 Tree based clustering.

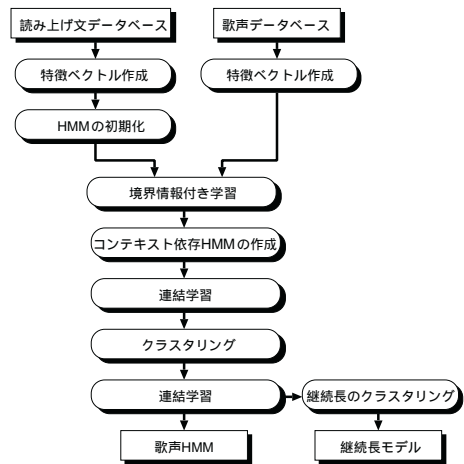


図 3 歌声システムの学習部  
Fig. 3 Training part of the system.

沿って木をたどることで、何らかのリーフノードに到達することができ、コンテキストの類似したモデルを選択することができる。

図 3 に、学習部の概要を示す。まず、テキスト読み上げ文による音声データベースから作成した不特定話者モデルから、歌声データのセグメンテーションを求める。次に歌声データベースを用いて楽譜情報を考慮した学習を行う。ここでは、歌声のモデル化に効果的なコンテキストとして以下に列挙するものを考えた歌声モデルを学習する。

- 歌詞：当該・先行・後続の音素名をコンテキストとし、母音・子音・有声音など音素に関する質問を適用。
- 音高：当該・先行・後続の音符の MIDI 音階値をコンテキストとし、音階の高低に関する質問を適用。
- 音長：当該・先行・後続の音符の長さ(100ms 単位で表したものを)をコンテキストとし、音符の長さに関する質問を適用。

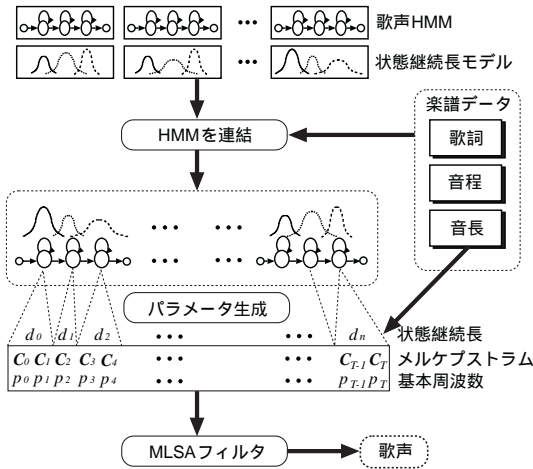


図4 歌声合成システムの合成部  
Fig. 4 Synthesis part of the system.

さらに、これらのコンテキストに基づいて MDL 基準を用いた決定木によるコンテキストクラスタリングを行い、あらゆるコンテキストの組合せに対応した歌声モデルを参照可能とする。一方、音素に対応する各歌声モデル内部の時間構造を表す状態継続長モデルは、HMM の各モデルの状態継続長を多次元ガウス分布でモデル化し、そのモデルパラメータは HMM の連結学習時に作られるトレリス上で求める<sup>13)</sup>。

2.2 合成部

図4に示す合成部では、楽譜情報(歌詞つき MIDI データ)を入力として歌声を合成する。まず、楽譜から得られる、歌詞、音高、音長情報に基づいて、歌声モデルから対応するモデルを選択する。次に、楽譜から与えられた各音符の長さを制約として、音符内の音素継続長および音素内部の状態継続長を、各モデルの状態継続長分布に基づいた尤度最大化基準により決定する。得られた状態系列から、次節で説明するパラメータ生成アルゴリズムによってメルケプストラムと基本周波数パラメータの列を生成する<sup>4)</sup>。最後に、生成されたパラメータに基づいて MLSA フィルタを励振させることで、歌声を合成する<sup>14)</sup>。

2.3 パラメータ生成アルゴリズム

連続出力分布型 HMM  $\lambda$  と状態遷移系列  $Q = \{q_1, q_2, \dots, q_T\}$  が与えられるとき、 $P(O|Q, \lambda)$  を最大にする音声パラメータ系列  $O = \{o_1, o_2, \dots, o_T\}$  を求めたい。ただし HMM の各状態は、状態  $q$  が  $d_q$  回継続する確率をガウス分布によりモデル化した状態継続長分布  $p_d(d_q)$  を持つものとする。また、簡単のため、HMM は単一出力分布型の left-to-right モデルを仮定している。

状態遷移系列  $Q$  が状態継続長分布から決定される場合、 $P(O|Q, \lambda)$  を最大化するパラメータ系列  $O$  はモデルの平均ベクトル系列と等しくなることは明らかであり、出力されるパラメータ系列は状態ごとに独立して決定されるため、各状態の境界において不連続が生じてしまう。

この問題を解決するため、静的特徴量と動的特徴量から構成される特徴パラメータ  $o_t = [c'_t, \Delta c'_t, \Delta^2 c'_t]'$  を導入する。なお、動的特徴量はデルタパラメータとも呼ばれ、音声認識で有効な特徴量であることが知られている。デルタパラメータは以下のように前後に隣接する静的特徴量  $c_t$  の線形結合によって表される。

$$\Delta^{(n)} c_t = \sum_{i=-L_-^{(n)}}^{L_+^{(n)}} w^{(n)}(i) c_{t+i}, \quad n = 1, 2 \quad (1)$$

このような制約の下で、 $P(O|\lambda)$  を最大化する静的なパラメータベクトル  $c_t$  からなる系列  $C = \{c_1, c_2, \dots, c_T\}$  は、線形方程式  $\partial \log P(O|Q, \lambda) / \partial C = 0_{TM}$  によって与えられ、これは文献4)に提案されている高速アルゴリズムによって逐次的に計算することができる。

このように生成されるパラメータベクトルの系列は、静的および動的特徴の統計量を反映したものとなる。

3. 歌声データベース

統計的アプローチによる音声のモデル化には、データベースが不可欠である。これまで、テキスト音声合成に利用可能な読み上げ文章などの音声データベースは様々なものが整備されているが、現在のところ歌声に関しては適切なものが入手できないため、新たに歌声データベースの整備を行った。

3.1 データの収録

童謡を中心とした 60 曲を用いて、男声 1 名の歌手手による歌声を収録した。単一指向性のコンデンサマイクを用い、スタンドに固定したマイク(ウインドスクリーンを装着)から口元までの距離を約 10 cm とし、歌手手は楽曲の MIDI 演奏をヘッドホンでモニタしながら、DAT へ録音した。

収録データの概要を表1に示す。なお、MIDI データについては、WWW 上で公開されているものを収集した。

収録された歌声は、DAT-Link+を用いて、サンプリング周波数 16 kHz、サンプルサイズ 16 bit、モノラル音声のデータとして計算機へ取り込んだ。なお、16 kHz ヘダウンサンプリングを行うことで聴覚上の

表 1 データベースの概要と収録機器

Table 1 Singing voice database and recording equipments.

歌い手	男性 1 名
楽曲	童謡など 60 曲 (合計で約 72 分)
DAT デッキ	SONY DTC-ZA5ES
マイクロホン	SONY C-355

劣化は避けられないが、現状では、分析合成系による音声品質劣化の影響のほうが大きいと考えられる。

### 3.2 データの整備

様々なコンテキストに基づいて歌声モデルを学習する場合、コンテキスト情報の信頼性は、合成音の品質に大きく影響を与えられられるため、テキスト音声合成においても、正確なコンテキスト情報の整備が重要な要素となっている。本手法では楽譜情報をコンテキスト要因の 1 つとして扱うが、実際の収録データには、歌詞の読み誤りや歌声と楽譜の音階が一致しないなどの誤りが含まれている。

予備的な実験から、誤ったコンテキスト情報を含むデータから学習を行った場合、部分的に音高を外した歌声が合成されるなどすることが分かっている。そこで、高精度な歌声のモデル化を行うため、歌声データベースの整備として以下の作業を行った。

- MIDI データの編集  
学習時に利用するため、主旋律データを作成するとともに、音符ごとに歌詞情報を付加する。
- 歌声に合わせた MIDI データの修正  
実際の歌声と伴奏に使用した MIDI データの間の、歌詞、音高の誤りについて、MIDI データを歌声に合わせて修正する。
- 音素境界ラベリング  
不特定話者の音素 HMM を用いて、歌声の音素境界の Viterbi アライメントを求め、音素境界を手作業で修正する。

## 4. 実験

前章で作成した歌声データベースを用いて歌声モデルを学習する。さらに、学習したモデルに対して学習データに含まれない楽曲を入力として、歌声を合成した。

### 4.1 学習データの作成

まず、収録した 60 曲の歌声についてメルケプストラム分析と基本周波数抽出を行い、HMM の学習データを作成した。基本周波数の抽出には TEMPO<sup>15)</sup> を用いた。メルケプストラム分析、基本周波数抽出に関する分析条件をそれぞれ表 2、表 3 に示す。

得られた分析データから、0~24 次のメルケプスト

表 2 データベースの分析条件 (メルケプストラム)

Table 2 Experimental condition for Mel-cepstral analysis.

学習データ	歌声
データ数	60 曲 (男声 1 名, 約 72 分)
サンプリング周波数	16 kHz
フレーム周期	5 ms
分析窓長	25 ms
窓関数	Blackman 窓
分析次数	24 次

表 3 データベースの分析条件 (基本周波数)

Table 3 Experimental condition for F0 analysis.

学習データ	歌声
サンプリング周波数	16 kHz
フレーム周期	5 ms
分析窓長	25 ms
上限/下限	370 Hz/70 Hz

ラム係数ベクトルと基本周波数値をフレームごとの静的特徴量とし、これに前後のフレームから計算される動的特徴量を加えたものを歌声モデルの学習データとした。t 番目のフレームのメルケプストラムの静的特徴をそれぞれ  $c_t$  としたとき、その動的特徴量  $\Delta c_t$  および 2 次動的特徴量  $\Delta^2 c_t$  は以下の式 (2), (3) から計算した。

$$\Delta c_t = \frac{1}{2}(-c_{t-1} + c_{t+1}) \quad (2)$$

$$\Delta^2 c_t = \frac{1}{4}(c_{t-1} - 2c_t + c_{t+1}) \quad (3)$$

基本周波数  $p_t$  についても同様に  $\Delta p_t, \Delta^2 p_t$  を求め、メルケプストラムと基本周波数の 2 つのストリームからなる学習ベクトルの次元数は合計 78 次元となる。

### 4.2 歌声モデルの学習

歌声データから抽出されたメルケプストラムと基本周波数を MSD-HMM によってモデル化する。HMM は単混合 5 状態の left-to-right モデルとし、音素はポーズと無音を含んだ 36 種類とした。

前節で作成した学習データを用いて、楽譜情報に依存したコンテキスト依存モデルを学習し、2.1 節で述べたとおり、歌詞から得られる音素のほか、MIDI データの音階表現値を利用した音高と、当該音素を含む音符の時間長を 100 ms 単位で分類した音長のコンテキストについて先行、当該、後続を考慮し、さらに MDL 基準に基づいたコンテキストクラスタリングを行い、各モデルの状態を共有化した。

なお、メルケプストラム、基本周波数、継続長の各モデルにコンテキストクラスタリングを行う際に、以下の 2 種類の手法を検討した。

手法 A: 各モデルで、2.1 節で述べたすべてのコンテ

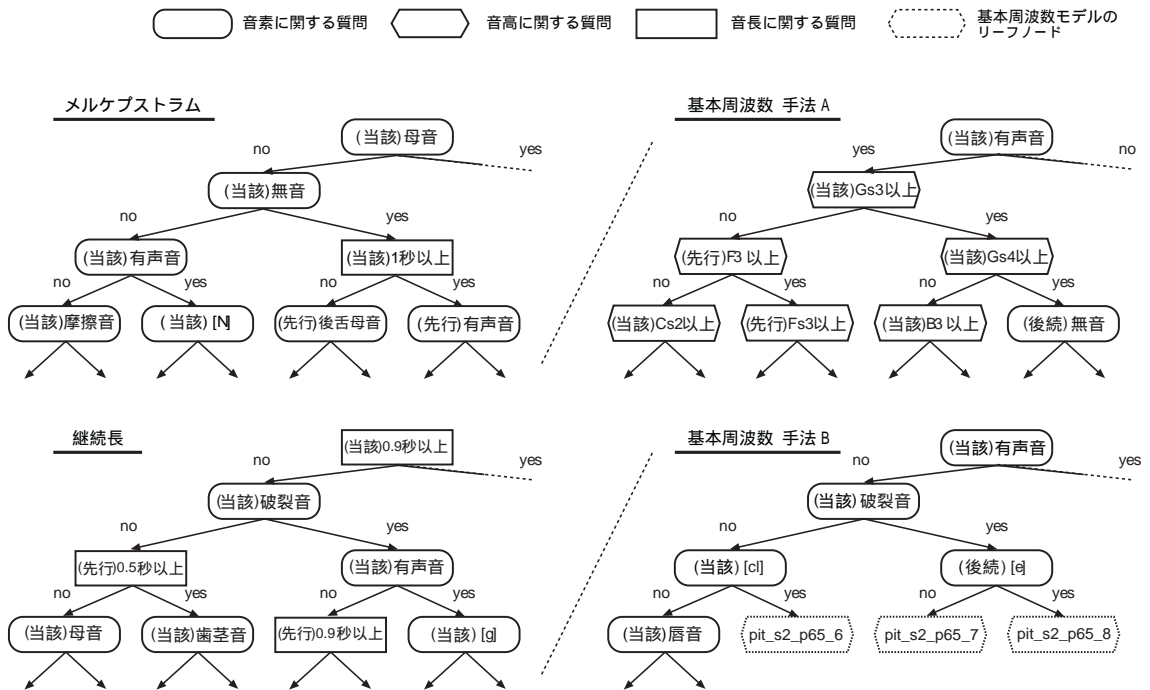


図 5 クラスタリングによって作成された各モデルの決定木

Fig. 5 Examples of decision trees.

キストを適用する。

手法 B: 基本周波数モデルに関してのみ, 当該音高別にクラスタリングを行う。

これは, 当該音高をコンテキストとした場合(手法 A), 異なる音高のデータが 1 つのクラスに分類され, 正しい音高が再現できなくなる可能性があるためである。

#### 4.3 歌声合成

学習データに含まれていない楽曲を用いて歌声を合成する。まず, 基本周波数のモデル化による自然性を確認するため, 以下の 2 つの手法から基本周波数系列を生成した。

手法 1: 楽譜から, 直接音階に相当する基本周波数系列を生成。

手法 2: 学習した基本周波数モデルから生成(提案法)。

なお, 手法 1 の音階に相当する基本周波数は, 以下の式 (4) から求める。p は MIDI 規格の音階を表す数値であり,  $p = 57$  が男声の基本周波数 110 Hz に相当する「ラ」の高さを表現している。

$$F_p = 110(2^{\frac{1}{12}})^{p-57} \quad (4)$$

#### 4.4 評価と考察

クラスタリングによって, メルケプストラム, 基本周波数(手法 A, B), 状態継続長のモデルから構築

された決定木の一部を図 5 に示す。主に, メルケプストラムでは音素に関する質問が多く適用されており, クラスタリング時に音高を固定しない場合(手法 A)の基本周波数モデルでは, 音高に関する質問が多く適用されていることが分かる。

また, 図 6 に生成された基本周波数(手法 1, 2)とスペクトル系列の一部を示す。手法 1 では音符単位に階段状に変化する平坦な基本周波数パターンが生成されているが, 手法 2 では HMM の内部状態にそって複雑に変化する基本周波数パターンが生成されている。また, 楽譜上の音階と比較して, 必ずしもそれに一致せず, 全体的にやや低くなっていることが分かる。これは学習データの歌い手が実際の楽譜よりも低く歌う傾向があり, その特徴がモデルに表れていることが考えられる。

#### 4.5 主観評価試験

合成された歌声の品質を評価するため, 手法 1, および手法 2 についてクラスタリング時の手法 A, B を考慮した手法, 計 3 通りの方法から, 学習データに含まれない 10 曲を合成した。各曲から 4 小節程度に分割した合計 32 のサンプルを切り出し, 被験者ごとにランダムに選択した 20 サンプルを用いて, 主観評価試験を行った。各被験者は, 各サンプルの自然性について 1~5 の 5 段階で評価を行った。10 名の被験者が

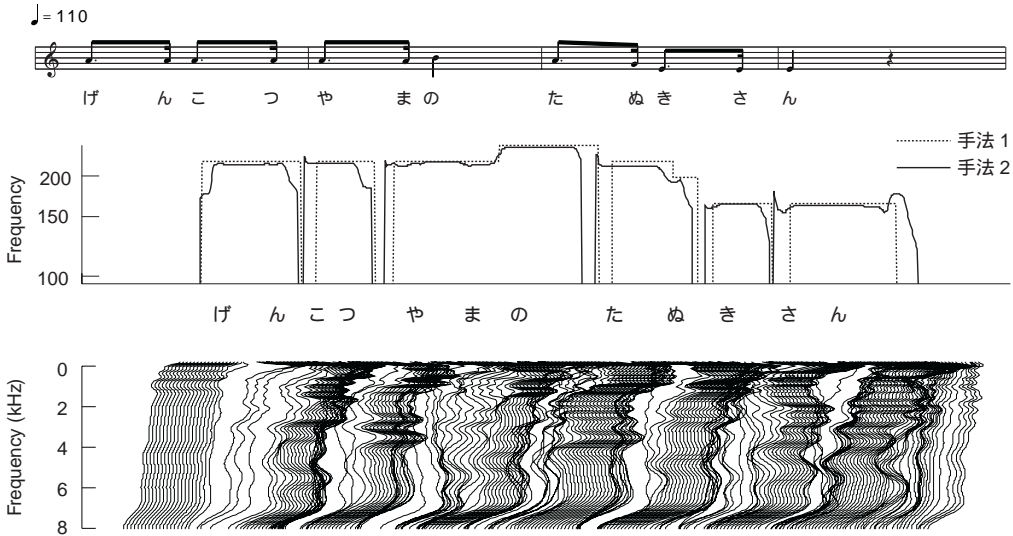


図 6 合成されたスペクトルと基本周波数パターン  
 Fig. 6 Example of generated spectra and F0 pattern.

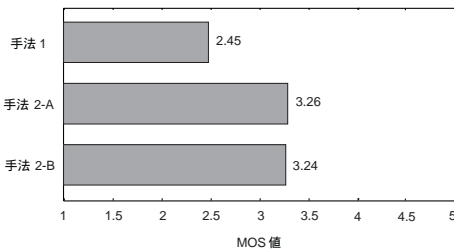


図 7 各手法の自然性についての評価  
 Fig. 7 Evaluation of naturalness for each methods.

ら得られた MOS 値を図 7 に示す。

試験の結果から、手法 2 の基本周波数モデルから生成した場合は、手法 1 の楽譜の旋律から生成した場合と比べ、高いスコアを得ている。これは図 6 の結果と同様に、HMM 内部の各状態ごとに、コンテキストに応じた変化を持つことが自然性の向上に大きく寄与していると考えられる。

また、基本周波数のクラスタリング手法による違い（手法 A, B による）に関しては、特に明確な差が得られていない。本来、音高を固定しないでクラスタリングを行う場合には、学習データが正しい音高のクラスに集まることは保証されないため、合成される歌声の音階がずれ、主観評価スコアが大きく劣化する可能性を持つことになるが、今回のサンプルでは、そのようなケースが見られなかった。これは、音高を固定しないクラスタリングにおいても、音高に関する質問が他のコンテキストより多く適用され、音高に応じたクラスターの分割が精度良く行えていることが原因の 1 つとして考えられる。

なお、非公式な結果として、本システムから合成された歌声を聴いた者の多くが、合成された歌声がデータベースの歌い手によるものと同定でき、また、当人の歌い方の特徴を感じ取ることができたことを付け加えておく。

### 5. まとめ

HMM に基づいた音声合成手法を拡張し、歌声合成システムを構築した。本システムでは、MSD-HMM によりスペクトルと基本周波数パターンを同時にモデル化することで、歌い手の特徴を再現した歌声合成が可能である。また、コンテキストクラスタリングにおいて、楽譜から得られる音高や音長をを利用することにより、歌声の精密なモデル化が可能であることを示した。

本論文では、新たに童謡 60 曲からなる歌声データベースの収録と整備を行い、それをを用いた実験から、学習データに含まれない楽曲に対しても、入力された任意の楽譜に対応した歌声合成が可能であることを示した。主観評価試験の結果からは、楽譜から規則的に生成した基本周波数パターンに対して、本手法で生成された基本周波数パターンの自然性の向上が確認できた。一方、非公式な結果として学習データ提供者の個性がよく再現されていることも確認することができた。

本手法は、ボコーダに基づいたシステムである点から、人間の歌声の代替となるような品質が達成できているわけではないが、たとえば、ロボットや玩具など

のエンタテインメントの分野においては、絶対的な品質よりは、むしろ、声質や歌い方などの個人性を再現することの方が重要であると考えられ、豊かな個性を実現する1つの方法として、幅広い応用が期待される。

また、本手法では精度の高いモデル化を行うため、歌詞、音高、音長をコンテキストとしたが、さらに発想記号や表情記号などをコンテキストとすることにより、より豊かな歌唱表現が可能と期待される。今回の実験では、曲によっては音符ごとの強弱が不自然になる問題が見られたが、音の強弱に関連したコンテキスト(フォルテ、ピアノなどの楽譜中の記号)を導入することにより、あわせてこの問題も解決されるものと考えられる。そのほか、長い音符において韻律の変動が平坦になる問題があったため、ピブラートなどの歌唱表現のモデル化に関する検討を行う必要がある。また、合成音に対して様々なフィルタ処理(たとえば文献16)を行うことでクリアな音質へと改善し、より実用性の高い歌声合成システムを開発することなどがあげられる。

謝辞 卒業研究を通して、本実験で使用したデータベースの収録、整備や実験にあたられた伊藤正典氏、石川ちさと氏の両名に感謝いたします。

### 参 考 文 献

- 1) 吉田由紀, 中嶋信弥: 歌声合成システム CyberSingers, 情報処理学会研究報告音声言語情報処理, Vol.25, No.8 (1995).
- 2) Macon, N.W., Jensen-Link, L.J., Oliverio, J. and Clements, M.A.: A Singing voice synthesis system based on sinusoidal modeling, *Proc. ICASSP*, Vol.1, pp.434-438 (1997).
- 3) 吉村貴克, 徳田恵一, 小林隆夫, 北村 正: HMMに基づく音声合成におけるスペクトル・ピッチ・継続長の同時モデル化, 信学論(D-II), Vol.J83-D-II, No.11, pp.2099-2107 (2000).
- 4) Tokuda, K., Yoshimura, T., Masuko, T., Kobayashi, T. and Kitamura, T.: Speech Parameter Generation Algorithm for HMM-based Speech, *Proc. ICASSP*, Vol.3, pp.1315-1518 (2000).
- 5) Tamura, M., Masuko, T., Tokuda, K. and Kobayashi, T.: Speaker Adaptation for HMM-based Speech Synthesis System Using MLLR, *Proc. 3rd ESCA/COCOSDA workshop on Speech Synthesis*, pp.273-276 (1998).
- 6) Yoshimura, T., Tokuda, K., Masuko, T. and Kobayashi, T.: Speaker Interpolation in HMM-based speech synthesis system, *Proc. EUROSPEECH*, Vol.5, pp.2523-2526 (1997).
- 7) Shichiri, K., Sawabe, A., Yoshimura, T., Tokuda, K. and Kitamura, T.: Eigenvoices for HMM-based speech synthesis, *Proc. ICSLP*, pp.1269-1272 (2002).
- 8) 徳田恵一, 益子貴史, 宮崎 昇, 小林隆夫: 多空間上の確率分布に基づいたHMM, 信学論(D-II), Vol.J79-D-II, No.7, pp.1579-1589 (2000).
- 9) 篠田浩一, 渡辺隆夫: 情報量基準を用いた状態クラスタリングによる音響モデルの作成, 信学技報, Vol.SP96-79, pp.9-16 (1996).
- 10) MIDI Manufactures Association. <http://www.midi.org/>
- 11) 徳田恵一, 小林隆夫, 斉藤博徳, 深田俊明, 今井聖: メルケプストラムをパラメータとする音声のスペクトル推定, 信学論(A), Vol.J74-A, No.8, pp.1240-1248 (1991).
- 12) Odell, J.J.: The use of context in large vocabulary speech recognition, Ph.D. Thesis, Cambridge University (1995).
- 13) 吉村貴克, 徳田恵一, 益子貴史, 小林隆夫, 北村正: HMMに基づく音声合成のための状態継続長モデルの構築, 信学技報, Vol.DSP98-85, No.262, pp.45-50 (1998).
- 14) 今井 聖, 住田一男, 古市千恵子: 音声合成のためのメル対数スペクトル近似(MLSA)フィルタ, 信学論(A), Vol.J66-A, No.2, pp.122-129 (1983).
- 15) Kawahara, H., Masuda, I. and Cheveigne, A.: Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based F0 extraction: possible role of a repetitive structure in sounds, *Speech Communication*, Vol.27, No.3-4, pp.187-207 (1999).
- 16) 吉村貴克, 徳田恵一, 益子貴史, 小林隆夫, 北村正: MMに基づく音声合成への混合励振源モデルとポストフィルタの導入, 信学技報, Vol.101, No.325, pp.17-22 (2001).

(平成 15 年 7 月 10 日受付)

(平成 16 年 1 月 6 日採録)



酒向 慎司(学生会員)

平成 11 年名古屋工業大学工学部知能情報システム学科卒業。現在同大学大学院博士後期課程在学中。視聴覚音声合成の研究に従事。電子情報通信学会, 日本音響学会各会員。





宮島千代美 (正会員)

平成 8 年名古屋工業大学工学部知能情報システム学科卒業。平成 13 年同大学大学院博士後期課程修了 (電気情報工学専攻)。同年名古屋工業大学知能情報システム学科助手。平成 15 年名古屋大学大学院情報科学研究科メディア科学専攻助手。工学博士。話者認識, マルチモーダル音声認識の研究に従事。第 17 回日本音響学会栗屋潔学術奨励賞受賞, 電子情報通信学会, 日本音響学会各会員。



北村 正 (正会員)

昭和 48 年名古屋工業大学工学部電子工学科卒業。昭和 53 年東京工業大学大学院博士課程修了。同年東京工業大学精密工学研究所助手。昭和 58 年名古屋工業大学工学部電子工学科講師。昭和 59 年同助教授。平成 7 年名古屋工業大学知能情報システム学科教授。工学博士。音声情報処理, マルチメディア情報処理の研究に従事。日本音響学会, IEEE, ISCA 各会員。



徳田 恵一 (正会員)

昭和 59 年名古屋工業大学工学部電子工学科卒業。平成元年東京工業大学大学院博士課程修了。同年東京工業大学電気電子工学科助手。平成 8 年名古屋工業大学知能情報システム学科助教授。工学博士。音声言語情報処理, マルチモーダル情報処理, 統計的学習理論の研究に従事。平成 13 年電気通信普及財団賞, 平成 13 年電子情報通信学会論文賞, 猪瀬賞各受賞。電子情報通信学会, 日本音響学会, 人工知能学会, IEEE, ISCA 各会員。

