

Guitarist Simulator : 演奏者の振舞いを 統計的に学習するジャムセッションシステム

浜 中 雅 俊^{†1,†2} 後 藤 真 孝^{†2,†3}
麻 生 英 樹^{†2} 大 津 展 之^{†2,†4,†5}

本論文では、ある人間の演奏者を模倣した仮想演奏者と、人間の演奏者とがインタラクティブに即興演奏できるジャムセッションシステムについて述べる。本研究の目的は、あたかも実在する人間の演奏者のような振舞いをする仮想演奏者を生成することである。従来の多くのジャムセッションシステムでは、仮想演奏者がどのような振舞いをするかは、アドホックなルール群などにより決定されていた。そのため、ルールのパラメータの調整により仮想演奏者に異なる個性を設定できても、実在する人間の演奏者と同じような振舞いをさせることは困難であった。本研究では、演奏者が聴取した演奏と現在弾いている演奏との関係、すなわち、演奏者の入出力関係を統計的に学習することにより、演奏者の振舞いのモデル（相手の演奏に対してどのような演奏をするかを定めるモデル）を獲得する手法を提案する。本手法の特長は、ジャムセッションの演奏記録さえあれば、それに参加した任意の演奏者の振舞いのモデルを獲得することができる点である。ギタートリオの MIDI 演奏を対象として、振舞いのモデルの学習機能を備えたジャムセッションシステム Guitarist Simulator を計算機上に実装した。実験の結果、MIDI ギターの演奏記録から、任意の 1 人の振舞いのモデルを学習可能なことを確認した。

Guitarist Simulator: A Jam Session System Statistically Learning Player's Reactions

MASATOSHI HAMANAKA,^{†1,†2} MASATAKA GOTO,^{†2,†3} HIDEKI ASOH^{†2}
and NOBUYUKI OTSU^{†2,†4,†5}

This paper describes a jam session system that enables a human player to interplay with virtual players, each of which imitates musical reactions of some human player. Our goal is to create a virtual player which reacts as if the actual human player does. Previous session systems have not been able to imitate such reactions, although those have parameters for altering a way of reacting. Our system can obtain the reaction model of a human player, namely, a characteristic way of reacting to the other players, by learning the relationship between MIDI data which the player listens to and MIDI data improvised by the player. It is not necessary to examine the target player directly: the session recording is all we need to build the model. The experimental results show that the proposed system called Guitarist Simulator dealing with a guitar trio can learn the reaction model of any participant from the MIDI recording of a session.

†1 日本学術振興会特別研究員 PD

Research Fellow of the Japan Society for the Promotion of Science

†2 独立行政法人産業技術総合研究所

National Institute of Advanced Industrial Science and Technology (AIST)

†3 科学技術振興事業団さきがけ研究 21「情報と知」領域

“Information and Human Activity,” PRESTO, Japan Science and Technology Corporation (JST)

†4 東京大学大学院情報理工学系研究科

Graduate School of Information Science and Technology, The University of Tokyo

†5 筑波大学大学院システム情報工学研究科

Graduate School, Doctoral Program Systems and Information Engineering, University of Tsukuba

1. はじめに

本研究の目的は、ある演奏者を模倣した仮想演奏者と人間の演奏者とが一緒にセッションすることができるシステムを実現することである。人間の演奏者はそれぞれ個性を持っており、相手の演奏に対してどのような演奏で反応するかは異なっている。ここでは、それを演奏者の振舞いと呼ぶ。あたかも実在する人間の演奏者のような振舞いをする仮想演奏者を生成するためには、演奏者個人の振舞いをモデル化する必要がある。演奏者の個性の違いは、このような振舞いの違いに表れると考えられるので、振舞いのモデル化が可能

となれば、親しい演奏者や、自分よりも演奏能力の高い演奏者、すでに亡くなった演奏者のモデルを用いた仮想演奏者といつでもインタラクションすることができるし、自分自身のモデルを用いた仮想演奏者とジャムセッションを行うことも可能となる。

従来のジャムセッションシステム^{1),2)}では、人間の演奏に追従した演奏を仮想演奏者にさせることに主眼がおかれていたため、仮想演奏者に個性を持たせるには至らなかった。一方、文献 3), 4) では、個性データベースと呼ばれるルール群を導入することにより個性の違いを設定することを可能としていた。文献 5), 6) では、システム外部から変更可能なパラメータを複数用意することにより、各演奏者が主導権を握る程度を様々に変化させることを可能としていた。しかし、これら^{3)~6)}は、パラメータやルールの調整を行うことにより、異なる振舞いの仮想演奏者を設定することはできても、実在する人間の演奏者の振舞いを模倣するようなモデルを設定することは困難であった。

演奏者の振舞いを模倣するためには、演奏者の入出力関係を学習により獲得する必要がある。文献 7), 8) では、演奏パターンの入出力関係をニューラルネットワークで学習することにより、人間の 8 小節の演奏に対して、8 小節の演奏パターンを出力するシステムを実現した。この研究は、演奏者のモデルを学習により獲得していた点が優れていたが、入出力関係を学習するうえで、どのような物理特徴量が重要であるかは検討されていなかった。また、人間とシステムが交互にソロを弾くという前提条件があったため、本研究のように両者が同時にソロや伴奏を弾くことはできなかった。さらに、モデルを学習する際にも、8 小節の入出力データをあらかじめ用意しなければならないという問題があった。実際のジャムセッションの演奏では、相手のどのフレーズに対して、どのフレーズで反応したかは明らかではなく、文献 7), 8) の手法を用いて、演奏記録から演奏者モデルを獲得することは困難であった。

これに対し本研究では、演奏記録から演奏者の入出力関係を統計的に学習することにより、演奏者の振舞いのモデルを獲得することを可能にする。このとき、この入出力関係を MIDI データなどから、直接統計的に求めようとすると、膨大な量のデータセットが必要となり、限られた長さの演奏記録からモデルを獲得することが困難である。そこで本研究では、入力と出力の物理的な特徴をそれぞれ一段抽象化した 2 つの主観空間（主観を表す空間）を導入する。その際、演奏記録を有効に用いることにより、どのような物理特徴量



図 1 人間の演奏者と仮想演奏者とが置換可能なセッションのモデル
Fig. 1 A session model in which either human or computer can be selected for each player.

が重要であるか検討した。そして、演奏者の振舞いのモデルを、入力側の主観空間から出力側の空間への関数として RBF ネットワークを用いて学習する。仮想演奏者は得られた関数を用いて、聴取した入力演奏から出力演奏を決定できる。

本研究では、文献 5) で後藤らが提唱しているセッションのモデルをさらに発展させ、人間の演奏者、仮想演奏者あわせて 3 人がギターでセッションするという形式で、人間の演奏者とその振舞いを模倣した仮想演奏者とが置換可能となるシステム *Guitarist Simulator* を実装した。3 人の人間の演奏者の MIDI 演奏記録をもとに学習を行ったところ、それぞれの演奏者の振舞いのモデルが獲得できた。また、得られたモデルを用いて仮想演奏者にセッションをさせたところ、ソロや伴奏を交代しながら、モデルに基づき実際の人間の演奏者に近い振舞いをすることが確認できた。

2. 演奏者を模倣するセッションシステム： Guitarist Simulator

本セッションシステムは、仮想演奏者が人間の演奏者と対等な立場ですべての演奏者の演奏を聴き、その演奏に反応するとともに自己主張することができるよう、文献 5) の全プレーヤが対等なセッションモデルに基づいて構成されている。本研究では、そのモデルをさらに発展させ、各演奏者が人間であるか仮想演奏者であるかが自由だけでなく、人間の演奏者に入れ替わってその振舞いを模倣した仮想演奏者に演奏させることが可能となるようにしている（図 1）。仮想演奏者は、自分を含めたすべての演奏者の演奏を聴き、その演奏に対して自分の振舞いを決定するため、相手

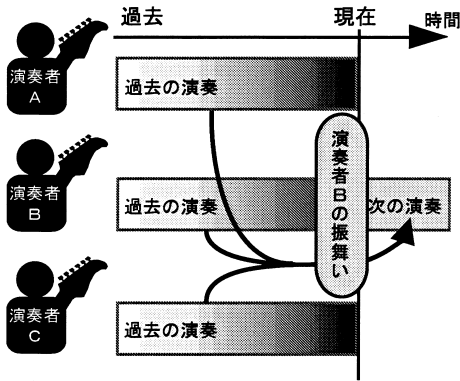


図 2 演奏者の振舞い
Fig. 2 Player's reactions.

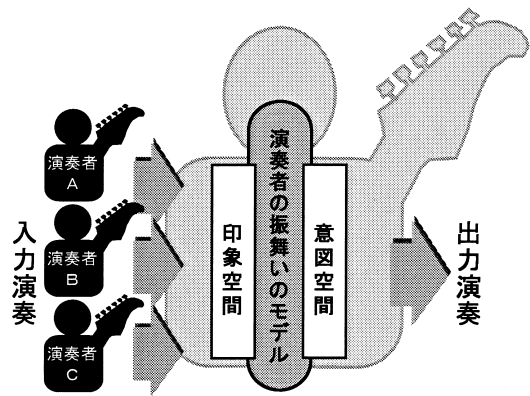


図 3 演奏者の振舞いのモデル
Fig. 3 Player's reaction model.

が人間であるか仮想演奏者であるかは問わない。

本研究では、以上のモデルをギタートリオのセッションに適用した。すべての演奏者にギターという同一の楽器を演奏させることにより、仮想演奏者の聴取過程でも複数の相手の演奏をまったく対等にとらえることができる。3人の演奏者は12小節1コーラスの典型的なブルース進行で、ソロや伴奏を繰り返しながら演奏する。調はA、拍子は4分の4拍子、テンポは120で一定、同じコード進行を12小節周期で繰り返すものとする。システムの入出力にはMIDIを用いる。

2.1 演奏者の振舞い

本研究では、全演奏者の過去の演奏から、次に自分がどのような演奏をするかを決める過程を「演奏者の振舞い」と定義する(図2)。システムでは、演奏の入出力にMIDIを用いているため、MIDI入力演奏からMIDI出力演奏を決定するものを演奏者の振舞いとして学習する。そのために、ここでは「印象空間」と「意図空間」という2つの主観的な空間を導入する。

システムが適切な出力を出すためには、セッション時のMIDI入力データと同じような入力が、学習時のMIDIデータに含まれていることが望ましい。しかし、実際のセッションの演奏では、MIDIデータのレベルで同じ演奏が再現される確率はきわめて低く、様々な入力に対応するためには、膨大な量のMIDI入出力データの組が必要となり、入出力関係を学習することが困難である。本研究では、この問題の解決法として入力を一段抽象化した、印象空間を導入する(図3)。

一方、実際のセッションの演奏では、同じような印象を与える演奏に対して、演奏者がMIDIデータのレベルでまったく同じ演奏で反応する可能性は非常に低く、複数の出力が考えられるため、入出力関係を学習することが困難である。そこで本研究では、この問題の解決法として出力を一段抽象化した意図空間を導入

する。

印象空間と意図空間の導入には、上記のような理由のほかに、抽象化により情報を圧縮して純化することで、一緒に演奏する演奏者の違いに依存しない演奏者の振舞いを獲得したいという狙いもある。

2.1.1 印象空間の導入

印象空間は、演奏から受ける印象を表すのに適していると思われる言葉(印象語と呼ぶ)を座標軸に持つ空間である。ジャムセッション中のある時刻において、ある1人の演奏者の演奏から受ける印象は、印象空間上の1点で表される。以下、これを印象ベクトルと呼ぶ。印象ベクトルは、入力演奏から得られる物理的な特徴量を座標軸に持つ空間(物理特徴空間)から印象空間への写像によって決定される。このとき、演奏の物理的な特徴も物理特徴空間上の1点で表されることになり、これを物理特徴ベクトルと呼ぶ。印象空間を導入することによって、MIDIデータ上では、異なる演奏であっても、それが同じ印象を与えるのであれば、同じ印象ベクトルとして表される。そのため、同じ印象を与える演奏がMIDIデータのレベルに比べて再現されやすくなり、セッション時の入力と同じような入力が学習データに含まれやすくなる。したがって、印象空間の導入により、限られた長さのMIDIデータから様々な入力に対応した入出力関係を学習することが可能となる。

2.1.2 意図空間の導入

意図空間は、仮想演奏者が演奏を生成するとき使用する1小節から8小節の長さの演奏パターンを、被験者実験の結果求めた主観的類似度(被験者の主観に基づく類似度)に基づき配置した空間で、どのような演奏を生成したいかを表す。ある1つの演奏パターンは意図空間上の1点で表される。以下、これを意図ベ

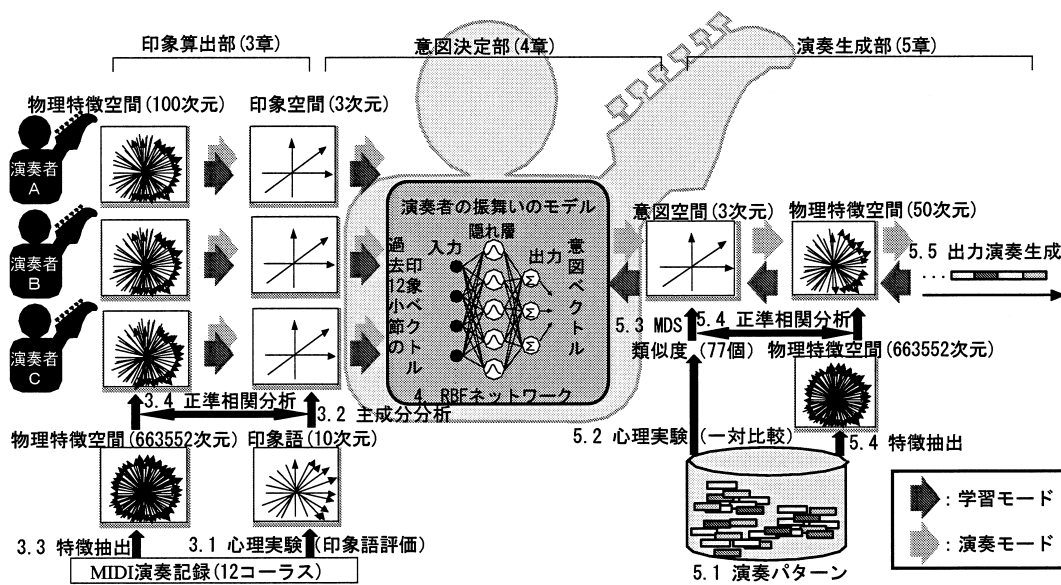


図 4 仮想演奏者の構成
Fig.4 Architecture of each virtual player.

クトルと呼ぶ。印象空間が人間の聴取過程に対応するのに対し、この意図空間は演奏生成過程に対応する。意図空間を導入することによって、演奏者の振舞いのモデルは、印象空間から意図空間への関数として学習することが可能となる。

2.2 仮想演奏者の構成

仮想演奏者は、1) 印象算出部、2) 意図決定部、3) 演奏生成部、の3つからなる(図4)。その詳しい内容については、3章、4章、5章で述べる。仮想演奏者は振舞いのモデルを学習するとき(学習モード)とそのモデルを用いて実際にセッションするとき(演奏モード)では、それぞれ異なった処理をする。学習モードは、非リアルタイムで実行され、印象ベクトルの時系列から意図ベクトルへの関数を3人の演奏者のジャムセッションの演奏記録から学習する。一方、演奏モードは、リアルタイムで実行され、学習によって得られたネットワークを用いて仮想演奏者の意図ベクトルを算出し、出力演奏を生成する。

2.2.1 学習モード

演奏者の振舞いのモデルを、人間の演奏者がセッションを行っている間に刻々と変化する印象ベクトルと意図ベクトルの組から統計的に学習する。具体的には、以下のような3段階の処理を行う(図5)。

まず、演奏者が聴取した演奏の MIDI データの時系列から物理特徴ベクトルを求め、それから印象ベクトルを算出する。印象ベクトルの算出は、セッションに参加しているすべての演奏者に対して別々に行うため、

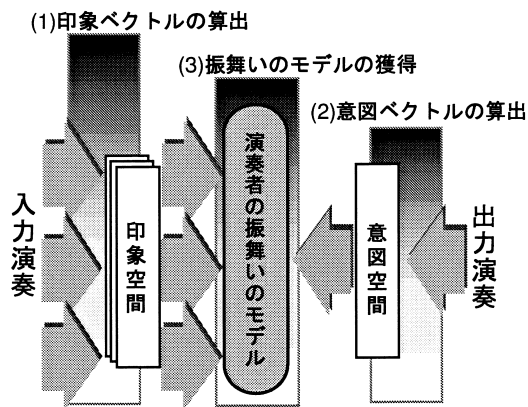


図 5 振舞いのモデルの学習の手順
Fig.5 Procedure of learning player's reaction model.

3つの印象ベクトルが得られる。次に、模倣しようとする演奏者が出力した演奏の MIDI データの時系列から物理特徴ベクトルを求め、その意図ベクトルを算出する。最後に、RBF ネットワークを用いて、演奏者の振舞いのモデルを、3つの印象ベクトルの時系列から意図ベクトルへの関数として獲得する。

2.2.2 演奏モード

演奏モードにおける仮想演奏者の動作は以下のとおりである。まず、3人の演奏者の入力演奏の物理特徴量を求め、得られた物理特徴量の時系列からそれぞれの印象ベクトルを算出する。次に、3人の印象ベクトルの時系列から、意図ベクトルを算出する。最後に、コード進行の制約を満たした演奏パターンの中から、

意図ベクトルに最も近いものを選択して出力する．演奏はシステムが出すカウント音によって開始する．

3. 印象算出部

印象算出部は、演奏者が聴取した演奏の MIDI データの時系列から物理特徴ベクトルを求め、それから印象ベクトルを算出する部分である．このような写像を行うためには、演奏者が「ある物理的特徴を持った演奏を聴いたとき、どのような印象を受けるか」を推定する必要がある．ここでは、演奏の物理特徴ベクトルと心理実験の結果得られた印象ベクトルとの間の相関関係を正準相関分析を用いて求めることにより、印象ベクトル算出のための写像を得ている．正準相関分析 (Canonical Correlation Analysis) とは、2 組の計測ベクトル間の相関関係を分析する場合に用いる手法で、複数の変数からなる 2 変数群 (ここでは印象ベクトルと物理特徴ベクトル) それぞれについて線形合成変数を求め、2 つの合成変数の相関 (正準相関) が最も大きくなるように重みをつけるものである．

3.1 印象ベクトルの獲得のための心理実験

被験者 (音楽経験者 1 人) に演奏を聴かせながら、1 小節の 1/12 の時間分解能で、演奏を表す 10 語 (安定感、異質感、開放感、緊張感、堅実感、重厚感、爽快感、存在感、平凡感、躍動感) の印象語がふさわしいかどうかを 7 段階で評定させる実験を行った．被験者は、1 小節の 1/12 の時間窓に含まれる音だけでなく、それ以前の音を含めてその時点での印象を評価する．用いた演奏は、ギタリストがテンポ 120 で行ったセッションの MIDI 演奏記録 (長さ 12 コーラス) の、ある 1 人分のパートの演奏である．10 語の印象語は、演奏者向きの音楽雑誌で使用頻度の高かった演奏を表す語の中から、特定の物理的特徴量との関わりが深いと考えられる語を選択している．このような印象語を選択することにより、評定の個人差をできる限り小さくすることを狙っている．

被験者は、実験用に作成したアプリケーションに表示された横軸の時間軸に沿って、印象が変化したと思われる点で縦軸に評定値をプロットしていく．印象の変化がなく、プロットしてない区間については、直前の評定値と同じ評定とする．演奏を 1 回聴いただけでは、すべての時刻における評定を入力することは難しいため、以上の操作は繰り返し行う．その際、前回までの評定結果は、画面上に表示される．

実験の結果、録音演奏から 1728 個 (12 コーラス × 12 小節 × 12 個) の印象ベクトルを得た．被験者がセッションの展開を記憶して評定を行うことを防ぐため、録音演奏を繰り返し聴く前に、ジャンルが異なる音楽 CD を 5 分以上聴かせた．しかし、それでもセッションの初めや終わりは、展開を記憶してしまいやすいので、録音演奏の前後 1 コーラス部分を除いた、1440 個 (10 コーラス × 12 小節 × 12 個) の印象ベクトルを用いることにした．実験の所要時間は、約 80 時間であった．

評定の時間分解能を 1 小節の 1/12 としているのは、被験者がこれより細かい時間分解能で印象の評定をすることは困難であると考えたためである．1 小節の 1/12 の時間分解能で評定した結果、印象ベクトルの連続的な変化が得られたため、その評定結果を用いることにした．演奏モードにおける印象ベクトルの算出の時間分解能も、評定の時間分解能と同様に、1 小節の 1/12 とした．

3.2 印象語の選択

演奏の物理特徴ベクトルと心理実験の結果得られた印象ベクトルとの間の相関関係を求める際、印象ベクトルはできる限り低次元であることが望ましい．印象ベクトルの低次元化は、演奏者の振舞いのモデルを効率良く学習するために必要である．ここでは、主成分分析を行い重要な印象語を選択することにより印象ベクトルを低次元化した．主成分分析 (Principal Component Analysis) とは、多変量の計測値から変量間の相関をなくし、しかもより低次元の変量によって元の計測値の特性を記述するための手法である．

印象ベクトルの主成分分析を行った結果、累積寄与率が 84% までの主成分が 3 つ得られ、第 1、第 2、第 3 主成分がそれぞれ、存在感、躍動感、重厚感に近接していた．したがって、この 3 つの印象語でほぼ印象空間全体を表現できることが分かった．以下、この 3 つの印象語を座標軸に持つ空間を印象空間とし、印象ベクトルも 3 次元で表すことにする．

3.3 様々な物理特徴量の抽出

正準相関分析する際、2 つの合成変数の相関をできる限り大きくするためには、印象ベクトルと相関の高い物理特徴量を見つける必要がある．しかし、どのような物理特徴が、演奏から受ける印象に大きく影響しているかは明らかでない．そこで本研究では、独立な様々な種類の物理特徴量を用意することにした．具体的には、以下のような処理を行った．

まず、過去 12 小節の演奏に含まれるすべての音に対し、1/48 小節の分解能でその時間窓に入った音の、

たとえば、緊張感という印象語はコードのテンションノート、安定感という印象語は、コードのルートノートとの関連が深い．

ペロシティ、音高、ピッチベンドの値を求める。ペロシティは、楽器を鳴らす速さ（鍵盤を押す速さ）で、0 から 127 の値で表される。音高は、音の高さを指定する MIDI のノート番号で、0 から 127 の値で表される。ピッチベンドは、音の高さを連続的に上下に変化させるためのパラメータで、-8192 から +8192 の値で表される。

次に、その値と、音量フィルタ、音域フィルタ、Note フィルタという 3 つのフィルタの出力を AND 結合した値の積を求める。音量フィルタ、音域フィルタはそれぞれ 3 つの関数からなり、Note フィルタは 512 個の関数からなる。そして各フィルタはそのうち 1 つの関数が選択されるため、3 種類のフィルタを AND 結合した値は、4,608 通り (3 (音量フィルタ) × 3 (音域フィルタ) × 512 (Note フィルタ)) 求められる。ここで AND 結合を採用した理由は 3.4 節で正準相関分析する際に、AND 結合、OR 結合、XOR 結合の 3 通りを試みた結果、AND 結合によって得られた物理特徴量には正準相関の値が高くなるものが多く含まれていたのに対して、OR 結合、XOR 結合により得られた物理特徴量には正準相関の値が高くなるものがほとんど含まれていなかったためである。

以下、それぞれのフィルタについて説明する。

(1) 音量フィルタ

音量フィルタは、以下の 3 つの関数からなり、音量の大小の情報を物理特徴量に取り込む意味がある。音量の小さな演奏と音量の大きな演奏がよく分離できるように閾値をペロシティ 91 に設定した。

- 音量の小さな演奏（主に伴奏）を抽出する関数
ペロシティが 90 以下の場合に 1 を出力し、そうでない場合に 0 を出力する。
- 音量の大きな演奏（主にソロ）を抽出する関数
ペロシティが 91 以上の場合に 1 を出力し、そうでない場合に 0 を出力する。
- 音量の小さな演奏と大きな演奏の両方を抽出する関数
つねに 1 を出力する。

(2) 音域フィルタ

音域フィルタは、以下の 3 つの関数からなり、音域の高低の情報を物理特徴量に取り込む意味がある。高音の演奏と低音の演奏がよく分離できるように閾値を音高 55 に設定した。

- 低音の演奏（主に伴奏）を抽出する関数

音高が 54 以下の場合に 1 を出力し、そうでない場合に 0 を出力する。

- 高音の演奏（主にソロ）を抽出する関数
音高が 55 以上の場合に 1 を出力し、そうでない場合に 0 を出力する。
- 低音の演奏と高音の演奏の両方を抽出する関数
つねに 1 を出力する。

(3) Note フィルタ

Note フィルタは、以下で述べる 9 個の関数の AND 結合をとった 512 種類 ($\sum_{n=0}^9 {}_9C_n$) であり、和音の情報を物理特徴量に取り込む意味がある。

- 3 和音を抽出する関数
コードの 3 和音（1 度、3 度、5 度）のいずれかの音の場合に 1 を出力し、そうでない場合に 0 を出力する。
- 6 度、7 度を抽出する関数
6 度または 7 度の場合に 1 を出力し、そうでない場合に 0 を出力する。
- テンションノートを抽出する関数（4 種類）
テンションノート（9 度、11 度、13 度、#9 度）の場合にそれぞれ 1 を出力し、そうでない場合に 0 を出力する。
- ブルーノートの 3 度、5 度を抽出する関数
ブルーノートの 3 度または 5 度の場合に 1 を出力し、そうでない場合に 0 を出力する。
- ブルーノートの 7 度を抽出する関数
ブルーノートの 7 度の場合に 1 を出力し、そうでない場合に 0 を出力する。
- その他の音を抽出する関数
上記の 8 種類の関数がいずれも 0 を出力するような音の場合に 1 を出力し、そうでない場合に 0 を出力する。

こうして得られた値の、過去 12 小節分の長さを 12 等分し、それぞれの、平均値、最大値、最小値、一次関数近似の傾きを算出し、それらを物理特徴量とした。物理特徴量は全部で 663,552 個 (12 小節 × 3 (ペロシティ、音高、ピッチベンド) × 3 (音量フィルタ) × 3 (音域フィルタ) × 512 (Note フィルタ) × 4 (平均値、最大値、最小値、傾き)) 生成された。

演奏の連続性を考えるうえでは、仮想演奏者が物理

ブルーノートとは、ブルーノート・ペンタトニック・スケールに含まれる音で、ブルーノートの 3 度、5 度、7 度は、それぞれ b 3 度、 b 5 度、 b 7 度である。

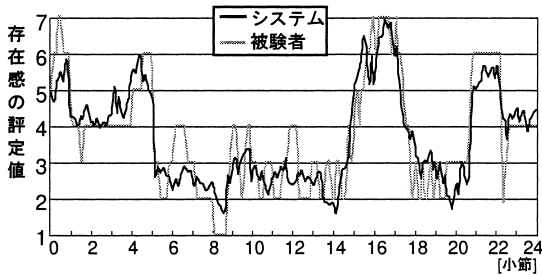


図 6 被験者とシステムが示した存在感の値の比較

Fig. 6 Comparing a 'presence' impression evaluated by the system with one evaluated by a subject.

特徴量を算出する時間分解能は、できる限り細かいほうが好ましいが、あまり分解能が細かいと物理特徴量算出の計算量が増大し、マシンパワーが不足するという問題が生じる。そこで本研究では、仮想演奏者のリアルタイム性とマシンパワーとの両者のバランスを考えて、1小節の1/48という時間分解能を設定した。

3.4 物理特徴量の取捨選択

正準相関分析を繰り返し行い、正準得点の大きさによって物理特徴量を取捨選択した。具体的には、以下のような処理を行った。物理特徴ベクトルと印象ベクトルとの正準相関分析を行うごとに、一番正準得点が低かった物理特徴量を1つ削除し、新たな物理特徴量を1つ加える操作を繰り返し、印象ベクトルと相関の高い物理特徴量を選択する。はじめに物理特徴ベクトルに与えられる物理特徴量は100個である。新たに追加する物理特徴量は、物理特徴ベクトルに残っている99個の物理特徴量のいずれに対しても相関が低いものとする。これは、相関の高い物理特徴量を加えても正準相関の値が変化しないだけでなく、相関行列の状態が悪くなり逆行列が求められなくなるなどの問題が起こるためである。物理特徴量間の相関は、単相関係数により求められる。単相関係数とは、2つの変数の共分散をそれらの標準偏差の積で割った値で、0のときに最も相関が低く、1のときに最も相関が高くなる。本研究では、物理特徴ベクトルに残っている99の物理特徴量との単相関係数がいずれも0.9以下のものを、新たな物理特徴量として追加することにした。正準相関分析を1,990,656回(663552×3回)行ったところ、正準相関の値が0.9を超えたため、そこで終了した。

3.5 印象ベクトルの獲得

印象ベクトルと物理特徴ベクトルの正準相関分析により、両者の写像が得られたため、新たな演奏が入力された場合でも、物理特徴ベクトルから印象ベクトルを求めることが可能となった。図6は、存在感につい

て、被験者の評定値とシステムの算出値とを、比較した結果である。92%の部分で両者の差は1以内であった。システムは、1小節に12回印象ベクトルを算出している。

4. 意図決定部

意図決定部は、演奏者の振舞いのモデルを学習し、得られたモデルを使って仮想演奏者の意図ベクトルを決定する部分である。この部分では、過去の印象ベクトルの時系列を入力し、各拍ごとに意図ベクトルを出力するようなネットワークを構成する。ソロや伴奏の切替えや演奏パターンの切替えは拍単位で行えば十分と考え、仮想演奏者が意図ベクトルを算出して次の演奏パターンに切り替えるのを、1拍単位で行うことにした。本研究では、4分の4拍子の演奏を扱っているため、意図ベクトルは1小節で4回算出される。

このようなことを行う学習アルゴリズムとして、本研究ではRBF(Radial Basis Function)ネットワークを採用した。RBFネットワークは、非線型関数を円形の等高線を持つ基底関数で展開する方法である。これを採用した理由は、データ間の補完能力に優れていること、ネットワークの重み算出が最小2乗法に帰着できて容易なことである。

本研究では、3人の演奏者から受け取る過去の印象ベクトルの時系列 x (L次元=3人×3次元×過去 n 小節分の印象ベクトル)を入力し、意図ベクトル $y=(\hat{y}_1, \hat{y}_2, \hat{y}_3)$ (3次元)を出力するようなRBFネットワークを構成する(図7)。意図ベクトルの次元数は、多次元尺度法により3次元としており、具体的には、5.3節で述べる。

学習用データは、12コーラス(1コーラスの長さは12小節)の演奏記録に含まれる、過去の印象ベクトル x とそれに対応する意図ベクトル y との対、 N 個($11=(12-1)$ コーラス×12小節×4拍=528個)である。したがって、入出力データの組は $(x^{(k)}, y^{(k)})(k=1, 2, \dots, N)$ となる。

ネットワークは M 個($M \leq N$)のユニットを持ち、各ユニットの中心 $c_i(i=1, 2, \dots, M)$ はL次元で表される空間上の1点を表す。各ユニットは、入力ベクトル $x^{(k)}$ とユニットの中心 c_i とのユークリッドノルム r を引数とするガウス関数 $\phi(r)$ を出力する(式(1))。

$$\phi(r) = \exp\left(-\frac{r^2}{\sigma^2}\right) \quad \sigma \neq 0, r \geq 0 \quad (1)$$

$$r = \|x^{(k)} - c_i\|$$

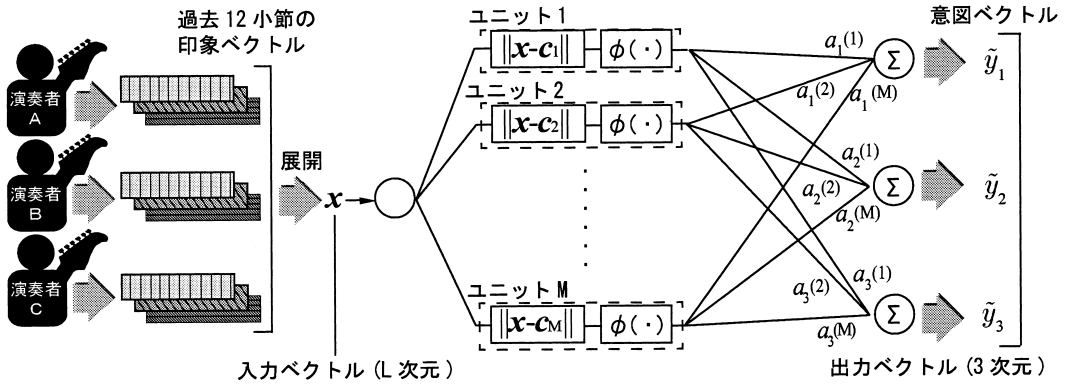


図 7 振舞いのモデルを学習する RBF ネットワーク

Fig. 7 RBF Network for learning player's reaction model.

ネットワークの出力 $\tilde{y}_j^{(k)}$ はそれに重み $a_j(i)$ を掛け加え合わせたものである (式 (2)).

$$\tilde{y}_j^{(k)} = \sum_{i=1}^M a_j(i) \phi(\|x^{(k)} - c_i\|) \quad (2)$$

入力ベクトル $x^{(k)}$ から任意の M 個をユニットの中心の候補 c_i として選択し、モデル化誤差を $e_j^{(k)}$ とするとネットワークの入出力関係は次式で与えられる.

$$y_j = P(x) a_j + e_j \quad (3)$$

ただし,

$$P(x) = \begin{bmatrix} p_1^{(1)} & p_2^{(1)} & \cdots & p_M^{(1)} \\ p_1^{(2)} & p_2^{(2)} & \cdots & p_M^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ p_1^{(N)} & p_2^{(N)} & \cdots & p_M^{(N)} \end{bmatrix}$$

$$p_i^{(k)} = \phi(\|x^{(k)} - c_i\|) \quad (k=1, 2, \dots, N, i=1, 2, \dots, M)$$

$$y_j = \begin{bmatrix} y_j^{(1)} \\ y_j^{(2)} \\ \vdots \\ y_j^{(N)} \end{bmatrix}, a_j = \begin{bmatrix} a_j(1) \\ a_j(2) \\ \vdots \\ a_j(M) \end{bmatrix}, e_j = \begin{bmatrix} e_j^{(1)} \\ e_j^{(2)} \\ \vdots \\ e_j^{(N)} \end{bmatrix}$$

4.1 ユニットの配置

入力ベクトル $x^{(k)}$ から出力 $y_j^{(k)}$ への寄与の大きいものを M 個求め、ユニットの中心 c_i として選択する. ユニットの配置は、 $y_1^{(k)}, y_2^{(k)}, y_3^{(k)}$ それぞれについて行う. このような処理を行うと、同じ入力データが複数含まれる場合に、重みが一意に定まらないという問題を回避できるだけでなく、計算時間の短縮および汎化性の向上にもつながる. しかし、寄与の大きいもの M 個を一度に求める方法はないため、本研究では、グラムシュミットの直交化法を用いて、寄与分が最大のものを 1 個ずつ選択する操作を M 回繰り返すことにより、ユニットの中心の選択を行った⁹⁾.

グラムシュミット法を使用すると P を直交ベクトル W , 上三角行列 B に分解することができる.

$$P = WB \quad (4)$$

ここで、 $W = [w_1, w_2, \dots, w_N]$ とすると、 w_i は、 P における第 i 直交成分であり、その寄与分は、 $(w_i^T y_j)^2 / (w_i^T w_i)$ で表される. 寄与分が最大の w_i を 1 つ選択し、残りの直交成分について再度グラムシュミット法を繰り返すことにより、出力に対して寄与の大きいユニットが選択される.

4.2 重みの学習

誤差 e_j を最小化するような a_j を最小 2 乗法で推定する. 式 (3) より 2 乗誤差は以下で与えられる.

$$\begin{aligned} \|e_j\|^2 &= \|y_j - P(x) a_j\|^2 \\ &= (y_j, y_j) - 2(y_j, P(x) a_j) \\ &\quad + (P(x) a_j, P(x) a_j) \end{aligned} \quad (5)$$

これを a_j の各成分で偏微分すると、

$$\begin{aligned} \frac{1}{2} \frac{\partial \|e_j\|^2}{\partial a_j} &= -P(x)^T y_j + P(x)^T P(x) a_j \\ &= 0 \end{aligned} \quad (6)$$

したがって、 a_j の推定値は式 (7) で与えられる.

$$a_j = (P(x)^T P(x))^{-1} P(x)^T y_j \quad (7)$$

5. 演奏生成部

演奏生成部は、あらかじめ意図空間上に配置してある演奏パターンから、仮想演奏者の意図ベクトルに合ったものを選択して出力する部分である. 意図空間は、演奏パターンの主観的な類似度評価に基づき構成した空間である. すなわち、ある演奏パターンは意図空間上の 1 点と対応し、似ていると感じられる演奏パターンは相互に近接し、そうでない演奏パターンは相互に離れるように配置した空間である. このような空

間を用いることにより、仮想演奏者の意図ベクトルにあまり変化がない場合には、似通った演奏を出力し続けることができるし、また仮想演奏者の意図ベクトルに一致した演奏パターンがデータベース上にない場合でも、その意図ベクトルに近い演奏を出力することができる。本研究では、このような意図空間を構成するために一対比較法、および多次元尺度法 (MDS) を用いる。そして、各出力演奏パターンの意図ベクトルとその物理特徴ベクトルとの写像は正準相関分析を用いて求める。

5.1 演奏パターンの生成

演奏パターンは、演奏者の演奏を MIDI 録音したもののから、8 小節の長さを上限として切り出して作成した。ソロのような演奏をしている区間では、演奏的にまとまったフレーズとなるように、長さを調整して切り取った。一方、伴奏のような演奏をしている区間では、1 コーラス内のあらゆる小節の頭から、区切りのよい位置までを切り取った。

5.2 一対比較実験

被験者 (3.1 節の被験者と同一人物) に、77 個の演奏パターンからランダムに選んだ 2 つの演奏パターンについて 1: 「似ていない」、2: 「少し似ていない」、3: 「少し似ている」、4: 「似ている」の 4 段階で評定させた。77 個のパターンの内訳は、1 コーラスのすべての小節から始めるソロのようなパターンを 5 種類ずつ計 60 個 (12 小節 × 5 種類)、A、E、D、それぞれのコードでの伴奏のようなパターンを 5 種類ずつ計 15 個 (3 コード × 5 種類)、2 種類のターンアラウンドである。評定の際、被験者には、あらかじめすべての演奏パターンを聴かせ、どのような演奏があるかを把握させておいた。

5.3 多次元尺度法

Kruskal の多次元尺度法¹⁰⁾ を用いて、一対比較実験により得られた類似度行列から意図空間を構成した。多次元尺度法は、演奏パターン j と演奏パターン k の類似度を δ_{jk} 、多次元空間での距離を d_{jk} としたとき、類似度の高い演奏フレーズほど距離が近くなるように多次元空間内の点の位置を決定する手法である (式 (8))。

$$\delta_{jk} > \delta_{lm} \quad \text{ならば} \quad d_{jk} \leq d_{lm} \quad (8)$$

このとき、式 (8) が成立しない度合いは、ストレス値 S で表される。意図空間は次元数とストレス値 S の両方が小さくなることが望ましい。そこで、1 次元から 5 次元の各次元数におけるストレス値を計算した。その結果、1、2 次元ではストレス値が高かったが、3 次元以降では、ストレス値が十分小さくなったため、

意図空間の次元数は 3 次元とした。

5.4 正準相関分析

正準相関分析を用いて、意図ベクトルと出力演奏の物理特徴ベクトルとの相関を求めた。用いた物理量は 3.3 節のものと同様である。分析の結果、正準相関の値が 0.88 となった。これにより、意図ベクトルから出力演奏の物理特徴ベクトルへの写像およびその逆写像が得られる。逆写像は、陽に求められ、演奏者の振る舞いのモデルを学習するとき、新たな演奏パターンをデータベース上に加えるときに用いる。

5.5 出力演奏の生成

出力演奏は、データベースにある 100 個の演奏パターンから選択されたものを切れ目なくつなげ合わせるにより生成される。仮想演奏者は、前の演奏パターンの終わりになると、コード進行の制約を満たした演奏パターンの中から意図ベクトルに最も近いものを次の演奏パターンとして選択する。データベースにある演奏パターンの半数はソロ演奏であり、仮想演奏者はどの小節からでも自分の意図に近いソロ演奏を開始することができる。

なお、演奏パターンの数が有限であることから、RBF ネットワークの出力する意図ベクトルと、実際に出力する演奏パターンの意図ベクトルとは厳密には一致しない。そこで、両者の距離を求め、RBF ネットワークの出力する意図ベクトルが、どれだけ実際に出力された演奏に反映されているかを調べた。3 軸の意図ベクトルのうち 1 成分について演奏中の両者の距離を求めた結果、平均 0.27 であった。これは、演奏中の意図ベクトルの変化の標準偏差が 7.91 であったのに対し十分小さい値である。

6. 実験結果

演奏者を模倣するセッションシステム Guitarist Simulator を計算機 (Pentium III 650 MHz, Windows98) 上に実装した。図 8 は、システムの出力画面である。画面は、横に 3 つのパネルが並んでおり、各パネルは 1 人の演奏者に対応している。パネルの一番上にあるラジオボタン (選択ボタン) は、そのパネルを仮想演奏者が使用するか、人間の演奏者が使用するかを切り替えるためにある。そして、ラジオボタンの下には、各パネル 2 つずつ立方体が描かれている。中央の薄い色の正方形が立方体の一番奥にある面であり、それを取り囲む 4 つの台形が立方体の上下左右の面である。上の立方体は印象空間を表し、その中にある球の位置が、印象ベクトルを表している。下の立方体は、意図空間を表し、その中にある球の位置が、意図ベク

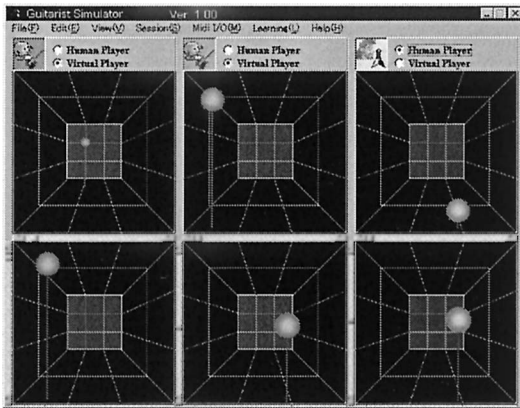


図 8 Guitarist Simulator の出力画面

Fig.8 Screen snapshot of Guitarist Simulator.

トルを表している。

6.1 振舞いのモデルの学習

3人の人間が、テンポ120のメトロノームの音にあわせて行ったセッションのMIDI演奏記録12コーラス分を用いて演奏者の振舞いのモデルの学習を行った。その際、RBFネットワークの入力次元数 L 、ユニット数 M 、式(1)の σ のそれぞれを適切に設定しないとモデル化誤差が大きくなり、汎化能力が低下する可能性がある。そこで、 L 、 M 、 σ をそれぞれ変化させモデル化誤差を計算した。その結果、 L が36次元(3人×4小節×3次元)、72次元(3人×8小節×3次元)、108次元(3人×12小節×3次元)の場合にはモデル化誤差が小さかったが、144次元(3人×16小節×3次元)以上ではモデル化誤差が大きかった。 L の次元数は大きいほど、過去の履歴をより反映した演奏ができるため、 L は108次元とした。 L を108次元とし、 M を1から200まで1刻み、 σ を0.5から70まで0.5刻みで、それぞれ変化させ、モデル化誤差が最も少なくなるパラメータを求めたところ、 $M=46$ 、 $\sigma=23.5$ が最適な値として求めた。他のサンプルデータについても、同様のパラメータの値で重みを計算できたため、演奏者の振舞いの学習を行う際には、つねにこの値を用いることにした。1回の学習は、ユニットの配置および重みの学習のみで済み、所要時間は5秒以内である。システムを用いて人間の演奏者がセッションを行うごとに、新たな演奏者の振舞いのモデルが獲得できた。

6.2 システムの評価

仮想演奏者A、Bと人間の演奏者Cでセッションを行い、仮想演奏者が人間の演奏者を模倣できているかを評価した。ここで、仮想演奏者A、Bは、人間の演奏者A、B、Cの3人で行った同じセッションの記

録から学習したもので、仮想演奏者Aは人間の演奏者Aの振舞いのモデルを持ち、仮想演奏者Bは人間の演奏者Cの振舞いのモデルを持つ。したがって、仮想演奏者Bが人間の演奏者Cのモデルを獲得できていれば、仮想演奏者Bと人間の演奏者Cはセッション中に同じような意図を持った演奏をしようとすると考えられる。

実験に参加したすべての人間の演奏者は、大学の音楽サークルの部員とOBで、5年以上のバンド経験を持つ。また、被験者にはシステムに関する説明はいっさいしていない。3人の演奏者の演奏は、それぞれ別々のスピーカから出力されており、誰がどの演奏を弾いたか分かるようになっている。人間の演奏者の音は生のギターの音を直接出力し、2人の仮想演奏者はそれぞれ異なるギターの音色をMIDI音源から出力した。

セッションを行った結果、人間の演奏者Cと仮想演奏者Bとは、ソロに入るタイミングが一致したり、類似したフレーズを演奏したりする傾向があった。12コーラス分の演奏の各時刻における人間の演奏者の意図ベクトルと仮想演奏者の意図ベクトルとの距離を合計し、その平均値を比較したところ、人間の演奏者Cと仮想演奏者Aでは、平均9.33と離れていたが、演奏者Cと仮想演奏者Bでは、平均3.56と近かった。以上より、仮想演奏者が人間の演奏者の振舞いをよく模倣できていることが確認された。参考までに、特に類似していると感じられた部分4コーラスでの3者の演奏意図の変化を、図9に例示する。

さらに、5人のギタリストにシステムとセッションした感想を求めたところ、すべてのギタリストが各仮想演奏者の振舞いの違いに気がついた。その中で、自分自身を模倣した仮想演奏者とジャムセッションした演奏者は2人いたが、ともに、仮想演奏者が自分の真似をしているようで演奏しにくいとの意見であった。

7. ま と め

本論文では、人間の演奏者の振舞いを模倣した仮想演奏者と人間の演奏者がインタラクションできるようなジャムセッションシステムについて述べた。本研究の主な意義は次の5点である。

- 統計的学習によりジャムセッションの演奏記録から演奏者モデルの獲得を実現した点。
- 心理実験の結果を用いて、印象空間と意図空間を構成し、振舞いのモデルの学習を可能とした点。
- 演奏の物理特徴量と、その演奏から得られる印象や意図を正準相関分析で関連づけるとともに、重要な物理特徴量の選択を行った点。

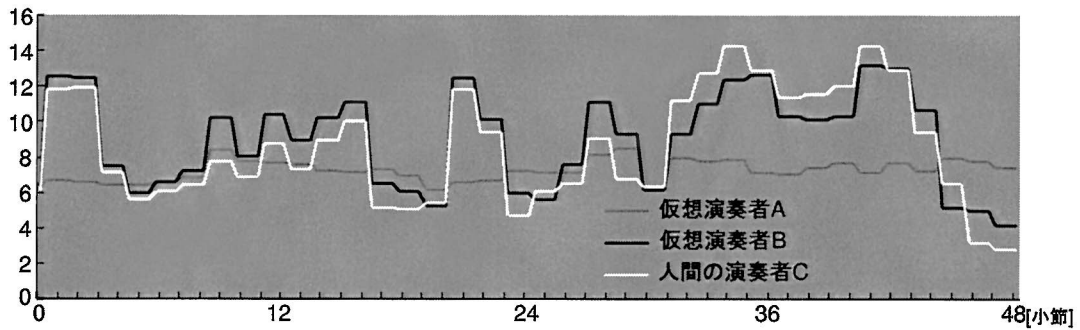


図9 3人の演奏者の意図ベクトルの変化(3軸の意図ベクトルのうち1成分のみを表示)

Fig.9 Transition of intention vectors of three players.

- 印象空間から意図空間への非線形関数を RBF ネットワークを用いて学習した点。
- 実際にリアルタイムで動作するシステムを実装し、実験により仮想演奏者が人間の演奏者の振舞いを模倣できていることを確認した点。

今後の課題として、第1に、学習アルゴリズムの改良があげられる。本システムでは、仮想演奏者の振舞いのモデルを過去12小節の印象ベクトルの履歴を用いて学習したため、セッション全体の展開や流れを把握した演奏は実現できていない。ジャズの曲などのように、コーラス数が決まっていたり、テーマを決めて演奏するような場合には、各演奏者がどのコーラスでどんな演奏をしたいかなど、今回扱ったものより、もっと大局的な演奏意図を考慮していく必要がある。今後、そのような大局的な演奏意図を扱うための手法について検討するとともに、大局的に見た時系列の変化も考慮にいたした学習アルゴリズムについて検討していく予定である。

第2の課題として、その他の個性の模倣があげられる。たとえば、本研究では、ある演奏とそれから受ける印象との関係は全演奏者で共通であると仮定して、1人の被験者の実験結果に基づいて印象ベクトルを算出している。この関係が演奏者によって異なるかを調査し、必要に応じて個性としてモデル化することは、今後の課題である。また、本研究では演奏の印象や意図のレベルでの演奏の入出力関係は、一緒に演奏する演奏者の違いに依存しないと仮定して振舞いのモデルを求めたが、演奏者間の相互作用など、演奏者の違いに依存する部分の個性のモデル化についても今後検討していく。他にも、グルーブ感と呼ばれるような演奏の心地よい揺らぎや、フレーズ、音色にも個性が表れると考えられる。今回は振舞いのモデル化に焦点を当てたが、今後、より詳細かつ総合的な個性の模倣に取り組んでいきたい。

謝辞 演奏者として実験に協力していただいた、齊田康公氏、千田真一氏、橋本大輔氏に感謝いたします。

参考文献

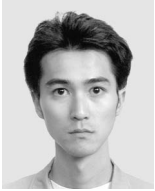
- 1) 青野裕司, 片寄晴弘, 井口征士: バンドライクな音楽アシスタントシステムについて, 情報処理学会研究報告, 94-MUS-8, Vol.94, No.103, pp.45-50 (1994).
- 2) 青野裕司, 片寄晴弘, 井口征士: アコースティック楽器を用いたセッションシステムの開発, 電子情報通信学会論文誌 D-II, Vol.J82-D-II, No.11, pp.1847-1856 (1999).
- 3) 和気早苗, 加藤博一, 才脇直樹, 井口征士: テンションパラメータを用いた協調型演奏システム—JASPER, 情報処理学会論文誌, Vol.35, No.7, pp.1469-1481 (1994).
- 4) 近藤欣也, 片寄晴弘, 井口征士: 音楽情報から奏者の意図を理解する伴奏システム JASPER++, 情報処理学会第46回全国大会, 1-373, 7Q-8 (1993).
- 5) 後藤真孝, 日高伊佐夫, 松本英明, 黒田洋介, 村岡洋一: 仮想ジャズセッションシステム: VirJa Session, 情報処理学会論文誌, Vol.40, No.4, pp.1910-1921 (1999).
- 6) 日高伊佐夫, 後藤真孝, 村岡洋一: すべてのプレイヤーが対等なジャズセッションシステム II, ベーシストとドラマーの実現, 情報処理学会研究報告, 96-MUS-14, Vol.96, No.19, pp.29-36 (1996).
- 7) Nishijima, M. and Kijima, Y.: Learning on Sense of Rhythm with a Neural Network—The NEURO DRUMMER, *Proc. ICMPC*, pp.78-80 (1989).
- 8) Nishijima, M. and Watanabe, K.: Interactive music composer based on neural networks, *Proc. ICMC*, pp.53-56 (1992).
- 9) Chen, S., Cowan, C.F.N. and Great, P.M.: Orthogonal Least Squares Learning Algorithm for Radial Basis Function Networks,

IEEE Trans. Neural Networks, Vol.4, pp.246–257 (1991).

- 10) Kruskal, J.B. and Wish, M.: *Multidimensional Scaling*, Sage Publications (1978).

(平成 15 年 7 月 9 日受付)

(平成 16 年 1 月 6 日採録)



浜中 雅俊 (正会員)

1998 年日本大学理工学部精密機械工学科卒業。2003 年筑波大学大学院工学研究科電子・情報工学専攻博士課程修了。現在、日本学術振興会特別研究員 PD, 独立行政法人産業技術総合研究所客員研究員として音楽情報処理の研究に従事。博士(工学)。2001 年情報処理学会山下記念研究賞, 2001 年 SCI(5th World Multiconference on Systemics, Cybernetics and Informatics) in Art 優秀論文賞, 2003 年筑波大学大学院優秀論文賞(博士課程長賞)各受賞。



後藤 真孝 (正会員)

1993 年早稲田大学理工学部電子通信学科卒業。1998 年同大学大学院理工学研究科博士後期課程修了。同年、電子技術総合研究所(2001 年に独立行政法人産業技術総合研究所に改組)に入所し、現在に至る。2000 年から 2003 年まで科学技術振興事業団さきがけ研究 21「情報と知」領域研究員を兼任。博士(工学)。音楽情報処理、音声言語情報処理等に興味を持つ。1992 年 jus 設立 10 周年記念 UNIX 国際シンポジウム論文賞, 1993 年 NICOGRAPH'93 CG 教育シンポジウム最優秀賞, 1997 年情報処理学会山下記念研究賞(音楽情報科学研究会), 1999 年平成 10 年電気関係学会関西支部連合大会奨励賞, 2000 年 WISS2000 論文賞・発表賞, 2001 年日本音響学会第 18 回粟屋潔学術奨励賞・第 5 回ポスター賞, 2002 年情報処理学会山下記念研究賞(音声言語情報処理研究会), 2002 年日本音楽知覚認知学会研究選奨, 2003 年インタラクシオン 2003 ベストペーパー賞各受賞。電子情報通信学会, 日本音響学会, 日本ソフトウェア科学会, 日本音楽知覚認知学会, ISCA 各会員。



麻生 英樹

1981 年東京大学工学部計数工学科数理工学コース卒業。1983 年同大学大学院工学系研究科情報工学専攻修士課程修了。同年電子技術総合研究所に入所。1993 年～1994 年ドイツ国立情報処理研究センター(GMD)客員研究員。現在、独立行政法人産業技術総合研究所情報処理研究部門主任研究員。知的学習システムに関する研究に従事。電子情報通信学会, 人工知能学会, 行動計量学会, 日本神経回路学会各会員。



大津 展之 (正会員)

1969 年東京大学工学部計数工学科卒業。1971 年同大学大学院数理工学専攻修士課程修了。同年電子技術総合研究所入所。以来、パターン認識, 画像処理, 多変量データ解析, 人工知能に関する数理的研究に従事。工学博士。数理工学情報研究室長, 首席研究官, 知能情報部長を経て, 現在。産業技術総合研究所フェロー。筑波大学連携大学院教授, 東京大学大学院情報理工学研究科教授を併任。電子情報通信学会, 日本行動計量学会等各会員。